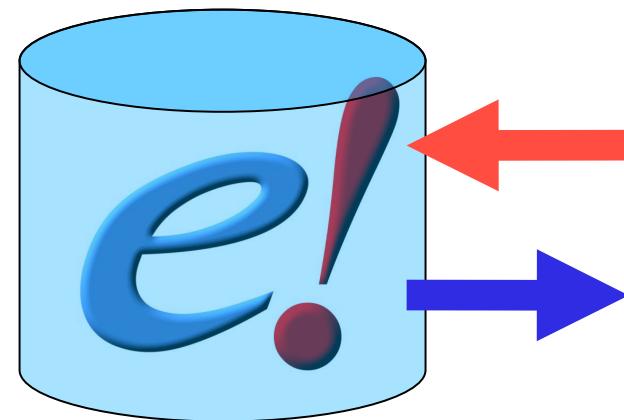




# Ensembl Core API



EMBL – European Bioinformatics Institute  
Wellcome Trust Genome Campus  
Hinxton, Cambridge, CB10 1SD, UK



# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references

# Ensembl API

- Written in Object-Oriented Perl.
- Used to retrieve data from and to store data in Ensembl databases.
- Foundation for the Ensembl Pipeline and Ensembl Web interface.



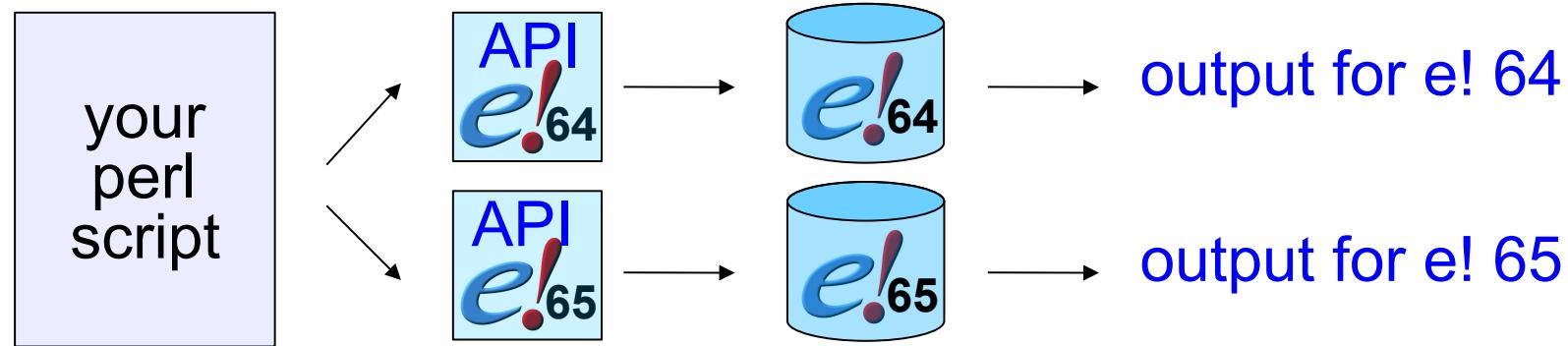
```
my $slice = $slice_adaptor->  
    fetch_by_region("chromosome", "13",  
                    31786617, 31872809);  
  
my $genes = $slice->get_all_Genes();  
  
foreach my $this_gene (@$genes) {  
  
    print $this_gene->stable_id, ":", "  
        $this_gene->start, " - ",  
        $this_gene->end, "\n";  
  
}  
}
```

# Why use an API?

- Uniform method of access to the data.
- Avoid writing the same thing twice: reusable in different systems.
- Reliable: lots of hard work, testing and optimisation already done.
- Insulates developers from underlying changes at a lower level (i.e. the database).

# Using the correct API version

API version **must** match database version. Old scripts using the API *should* continue working with a newer API.



Run script `ensembl/misc-scripts/ping_ensembl.pl` to test if you can contact the Ensembl database server. The script will help you resolve issues with your setup.

# An Alternative - REST

 Ensembl Home User Guide About the Ensembl Project Contact Ensembl

## Ensembl REST API Endpoints

### Comparative Genomics

Resource	Description
<a href="#">GET genetree/id/:id</a>	Retrieves Gene Tree dumps for a given Gene Tree stable identifier
<a href="#">GET homology/id/:id</a>	Retrieves homology information by ensembl gene id
<a href="#">GET homology/symbol/:species/:symbol</a>	Retrieves homology information by symbol

### Cross References

Resource	Description
<a href="#">GET xrefs/id/:id</a>	Perform lookups of Ensembl Identifiers and retrieve their external cross references in other databases
<a href="#">GET xrefs/name/:species/:name</a>	Performs a lookup based upon the primary accession or display label of an external reference and returning the information we hold about the entry
<a href="#">GET xrefs/symbol/:species/:symbol</a>	Looks up an external symbol and returns all Ensembl objects linked to it. This can be a display name for a gene/transcript/translation, a synonym or an externally linked reference. If a Gene's transcript is linked to the supplied symbol the service will return both Gene and Transcript (it supports transient links).

### Information

Resource	Description
<a href="#">GET assembly/info/:species</a>	Returns information about the current available assemblies in this given species
<a href="#">GET assembly/info/:species/:region_name</a>	Returns information about the given toplevel sequence region given to this endpoint

**<http://beta.rest.ensembl.org>**

# Installing the Perl API from FTP

```
# cd to a location to install Ensembl to
mkdir src
cd src

# Get the latest API from FTP (always the live version) and BioPerl 1.2.3
wget ftp://ftp.ensembl.org/pub/ensembl-api.tar.gz
wget http://bioperl.org/DIST/old_releases/bioperl-1.2.3.tar.gz

# untar both
tar zxvf ensembl-api.tar.gz
tar zxvf bioperl-1.2.3.tar.gz

# open up .bashrc or .profile and add the following
PERL5LIB=${PERL5LIB}:${HOME}/src/bioperl-1.2.3
PERL5LIB=${PERL5LIB}:${HOME}/src/ensembl/modules
PERL5LIB=${PERL5LIB}:${HOME}/src/ensembl-compara/modules
PERL5LIB=${PERL5LIB}:${HOME}/src/ensembl-variation/modules
PERL5LIB=${PERL5LIB}:${HOME}/src/ensembl-functgenomics/modules
export PERL5LIB

# Checking your installation
perl ${HOME}/src/ensembl/misc-scripts/ping_ensembl.pl
```

# Alternative Methods of Installation

- Ensembl tarball from <http://cvs.sanger.ac.uk>
  - Watch out for & character in Linux shells
- CVS checkout from cvs.sanger.ac.uk
- Ensembl Virtual Machine
  - [ftp://ftp.ensembl.org/pub/current\\_virtual\\_machine](ftp://ftp.ensembl.org/pub/current_virtual_machine)
  - 64bit only in OVA format

# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references

# Ensembl API – Object Types

We have two main object types in Ensembl API:

## 1. Data Objects

Talk to a particular row in a data table, such as the BRCA2 gene to get (or set) information related to this gene only.

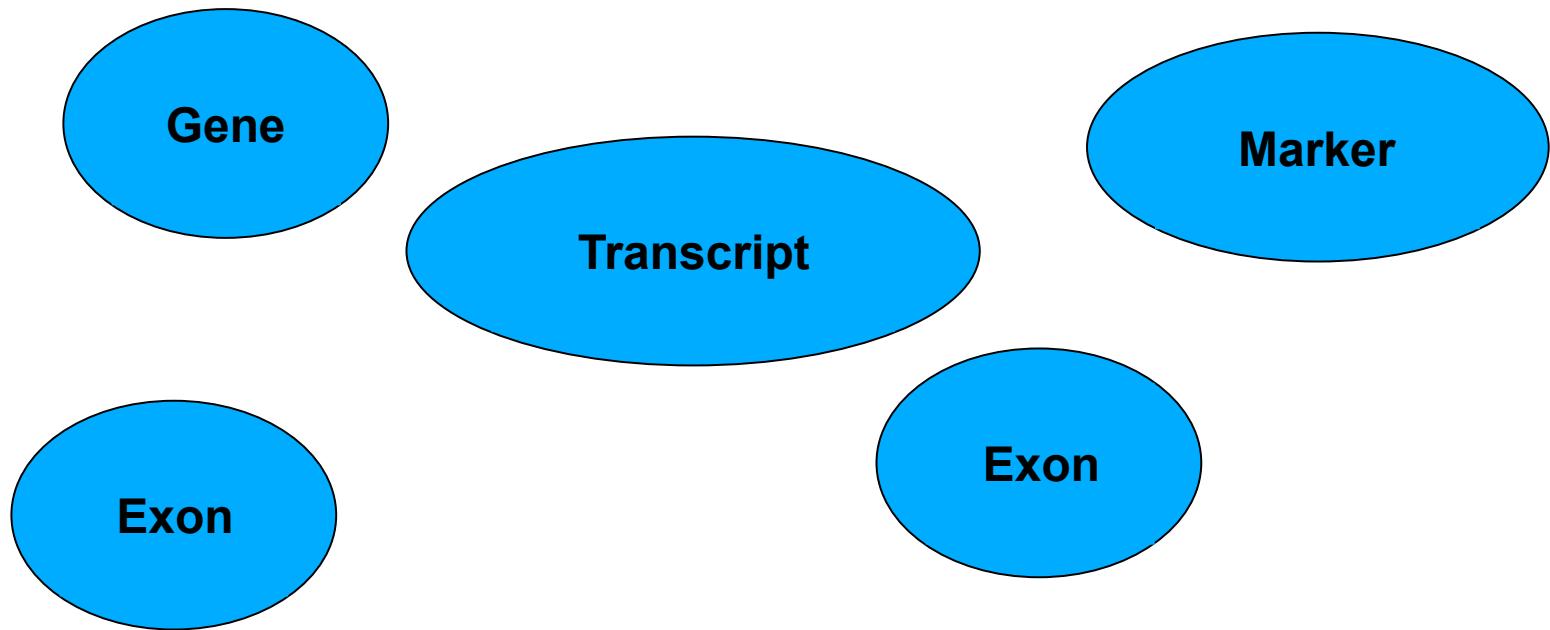
## 2. Object Adaptors

Talk to a particular data table, such as the gene table to retrieve or store genes in the gene table.

# Data Objects

Data objects model biological entities, e.g. genes, transcripts, exons...

A *Data Object* represents a piece of data that is (or can be) stored in the database.



# Object Adaptors

*Data Objects* are retrieved from and stored in the database using **Object Adaptors**.

Each *Object Adaptor* is responsible for creating objects of only one particular type. For instance:

- The **GeneAdaptor** is used to fetch **Gene** objects
- The **ExonAdaptor** is used to fetch **Exon** objects

*Object Adaptor* *fetch*, *store*, and *remove* methods are used to retrieve, save, and delete information in the database.

Two types of methods:

- **fetch\_by\_....** returns 1 object (or undef)
- **fetch\_all\_by\_....** returns a ref. to an array of objects (or ref. to an empty array)

# Object methods in Ensembl API

An Object has attributes and methods.

We avoid accessing object attributes directly, we use methods instead.

Methods are called using the “arrow” (->) operator:

```
my $exons = $gene->get_all_Exons();
```

Many methods can be used to either get or set an attribute value:

GET (no arg.):

```
my $gene_id = $gene->stable_id();
```

SET (new value):

```
$gene->stable_id("ENSG000000123152");
```

# Data Objects & Object Adaptors

```
# fetch a gene by its stable identifier using a gene
adaptor

my $gene =
    $gene_adaptor->fetch_by_stable_id('ENSG00000139618');

# print out the name of a gene and its stable identifier

print $gene->external_name(), "\n";
print $gene->stable_id, "\n";
```

# Ensembl Core modules

## Name space for the modules:

- Object modules start with **Bio::EnsEMBL**
  - Bio::EnsEMBL::Gene for gene objects
  - Bio::EnsEMBL::Exon for exon objects
- ObjectAdaptors start with **Bio::EnsEMBL::DBSQL**
  - the GeneAdaptor is a Bio::EnsEMBL::DBSQL::GeneAdaptor
  - The ExonAdaptor is a Bio::EnsEMBL::DBSQL::ExonAdaptor

## Naming conventions for the methods:

- For retrieving objects from the ObjectAdaptors:
  - fetch\_by\_ and fetch\_all\_by\_
- For retrieving Objects attributes:
  - get\_Object and get\_all\_Objects

# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. **Ensembl documentation**
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references

# Ensembl API Documentation (1)

Go to <http://www.ensembl.org/info/docs/Doxygen/index.html>

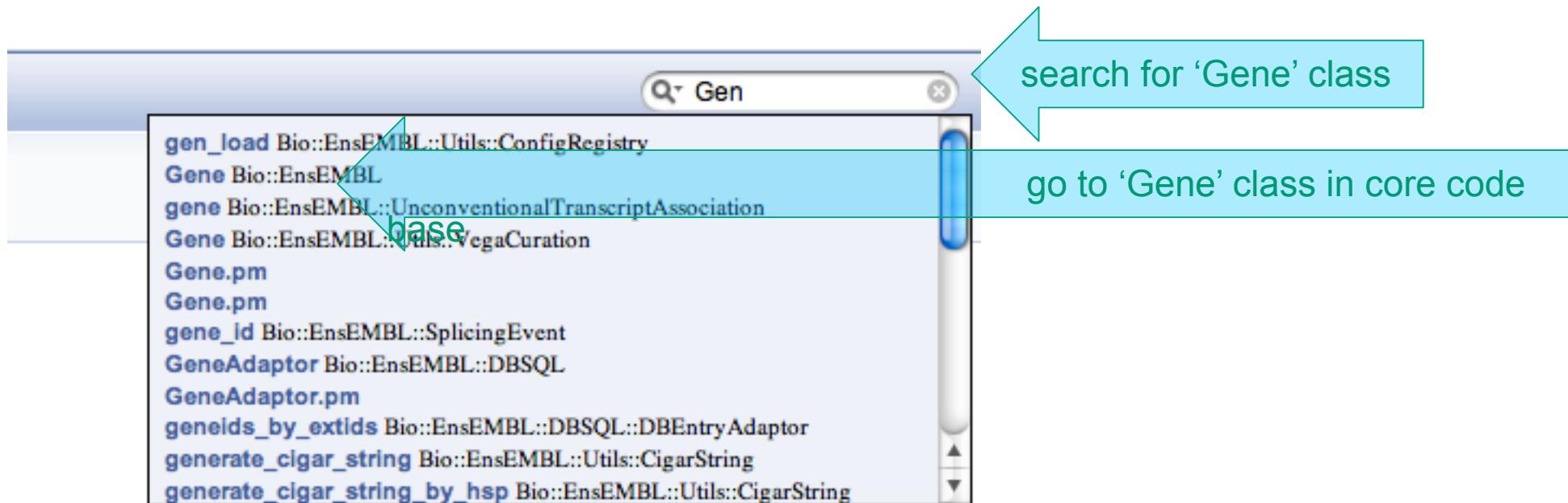
The screenshot shows the Ensembl website's navigation bar with links for BLAST/BLAT, BioMart, Tools, Downloads, and Help & Documentation. Below the navigation bar, a breadcrumb trail shows the user has navigated from the homepage to Help & Documentation and then to Doxygen Perl documentation. The main content area is titled "Doxygen perl module documentation". A teal arrow points from the text "Core code base" to the first item in the list of modules.

These are the components of the Ensembl code base. Click on a module below for API documentation:

- [Ensembl](#) - Ensembl core database API
- [Ensembl-analysis](#)
- [Ensembl-hive](#)
- [Ensembl-compara](#) - Ensembl comparative genomics API
- [Ensembl-external](#)
- [Ensembl-functgenomics](#) - Ensembl functional genomics API
- [Ensembl-pipeline](#) - Ensembl gene build pipeline\*
- [Ensembl-variation](#) - Ensembl variation data API
- [BioPerl](#)

# Ensembl API Documentation (2)

Search for classes or methods



# Ensembl API Documentation (3)

Breadcrumbs at the top of the page help you keep track of which code base you're in.

The screenshot shows a web-based API documentation interface. At the top, a breadcrumb navigation bar is visible with links: Home > Help & Documentation > Doxygen Perl documentation > Ensembl > Bio > EnsEMBL > Gene. Below this is a horizontal menu bar with tabs: Overview, Classes (which is selected and highlighted in dark blue), Namespaces, Files, Directory Browser, and Related Pages. Under the Classes tab, there are sub-links: Class Hierarchy, Class Index, Full Class List, and Methods from all packages. The main content area is titled "Bio::EnsEMBL::Gene Class Reference". It includes sections for "Inheritance diagram for Bio::EnsEMBL::Gene:", "List of all members.", "Class Summary", and "Synopsis". The Synopsis section contains the following Perl code:

```
my $gene = Bio::EnsEMBL::Gene->new(
    -START  => 123,
    -END    => 1045,
    -STRAND => 1,
    -SLICE  => $slice
);

# print gene information
print("gene start:end:strand is "
    . join( ":", map { $gene->$_ } qw(start end strand) )
    . "\n");
```

# Ensembl API Documentation (4)

Classes can *inherit* attributes and methods from other classes. For instance a Gene is also a Feature as it can be located on a region of DNA so it inherits from the **Bio::EnsEMBL::Feature** object.

The screenshot shows a web-based API documentation interface for the Ensembl BioPerl API. The top navigation bar includes links for Help & Documentation, Doxygen Perl documentation, Ensembl, Bio, EnsEMBL, and Gene. Below this is a menu bar with tabs for Overview, Classes (which is selected), Namespaces, Files, Directory Browser, and Related Pages. Under the Classes tab, there are links for Class Hierarchy, Class Index, Full Class List, and Methods from all packages. The main content area is titled "Bio::EnsEMBL::Gene Class Reference". It contains sections for Inheritance diagram, List of all members, Class Summary, and Synopsis. The Synopsis section displays Perl code for creating a new Bio::EnsEMBL::Gene object and printing its information. A large blue arrow points from the "Inheritance diagram for Bio::EnsEMBL::Gene" link to a list of inheritance sources on the left side of the page, which includes classes like Bio::EnsEMBL::Util::Eprof, Bio::EnsEMBL::Util::EprofStack, Bio::EnsEMBL::Util::Exception, Bio::EnsEMBL::Exon, Bio::EnsEMBL::DBSQL::ExonAdaptor, Bio::EnsEMBL::IdMapping::ExonScoreBuilder, Bio::EnsEMBL::External::ExternalFeatureAdaptor, Bio::EnsEMBL::DB::ExternalFeatureFactor, Bio::EnsEMBL::Feature, Bio::EnsEMBL::FeaturePair, Bio::EnsEMBL::DBFile::FileAdaptor, Bio::EnsEMBL::Mapper::Gap, and Bio::EnsEMBL::Gene. The "Bio::EnsEMBL::Gene" class is highlighted with a dark blue background.

expand to view classes which 'Gene' inherits from

```
my $gene = Bio::EnsEMBL::Gene->new( -START => 123, -END => 1045, -STRAND => 1, -SLICE => $slice ); # print gene information print("gene start:end:strand is " . join( ":", map { $gene->$_ } qw(start end strand) ) . "\n");
```

# Ensembl API Documentation (5)

Scroll down to see a list of available methods in the 'Gene' class. Can you find a method to retrieve all transcripts for a gene? What variable type does this method return?

[Help & Documentation](#) > [Doxygen Perl documentation](#) > [Ensembl](#) > [Bio](#) > [EnsEMBL](#) > [Gene](#)

Overview	Classes	Namespaces	Files	Directory Browser	Related Pages
<a href="#">Class Hierarchy</a>	<a href="#">Class Index</a>	<a href="#">Full Class List</a>		<a href="#">Methods from all packages</a>	

Bio::EnsEMBL::Mapping::Entry  
Bio::EnsEMBL::Utils::Eprof  
Bio::EnsEMBL::Utils::EprofStack  
Bio::EnsEMBL::Utils::Exception  
Bio::EnsEMBL::Exon  
Bio::EnsEMBL::DBSQL::ExonAdaptor  
Bio::EnsEMBL::IdMapping::ExonScoreBuilder  
Bio::EnsEMBL::External::ExternalFeatureAdaptor  
Bio::EnsEMBL::DB::ExternalFeatureFactor  
Bio::EnsEMBL::Feature  
Bio::EnsEMBL::FeaturePair  
Bio::EnsEMBL::DBFile::FileAdaptor  
Bio::EnsEMBL::Mapper::Gap  
**Bio::EnsEMBL::Gene**  
Bio::EnsEMBL::Utils::VegaCuration::GeneAdaptor  
Bio::EnsEMBL::DBSQL::GeneAdaptor

**Available Methods**

protected	<a href="#">_deprecated_transform ()</a>
public Bio::EnsEMBL::DBSQL::BaseAdaptor	<a href="#">adaptor ()</a>
public void	<a href="#">add_Attributes ()</a>
public void	<a href="#">add_DBEntry ()</a>
public DEPRECATED	<a href="#">add_DBLink ()</a>
public	<a href="#">add_sub_SeqFeature ()</a>
public void	<a href="#">add_Transcript ()</a>
public void	<a href="#">add_unconventional_transcript_association ()</a>
public Bio::EnsEMBL::Analysis	<a href="#">analysis ()</a>
public String	<a href="#">biotype ()</a>
public String	<a href="#">canonical_annotation ()</a>
public Bio::EnsEMBL::Transcript	<a href="#">canonical_transcript ()</a>
public	<a href="#">chr_name ()</a>
public	<a href="#">confidence ()</a>
public	<a href="#">contig ()</a>
public String	<a href="#">coord_system_name ()</a>
public String	<a href="#">created_date ()</a>

# Ensembl API Documentation (6)

Clicking on a method takes you to a description with associated arguments, return types and exceptions.

**public Bio::EnsEMBL::Slice Bio::EnsEMBL::Feature::slice ( )**

Arg [1] : (optional) **Bio::EnsEMBL::Slice** \$slice  
Example :

```
$seqname = $feature->slice()->name();
```

Description: Getter/Setter for the **Slice** that is associated with this feature. The slice represents the underlying sequence that this feature is on. Note that this method call is analogous to the old **SeqFeature** methods **contig()**, **entire\_seq()**, **attach\_seq()**, etc.

Returntype : **Bio::EnsEMBL::Slice**

Exceptions : thrown if an invalid argument is passed

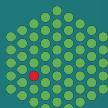
Caller : general

Status : Stable

## ▶Code:

[click to view](#)

Reimplemented in **Bio::EnsEMBL::Exon**, and **Bio::EnsEMBL::Map::DitagFeature**.



# Ensembl API Documentation (7)

Clicking on **Code**: shows the method's implementation

## ▼Code:

```
sub slice {
    my ( $self, $slice ) = @_;
    if ( defined($slice) ) {
        if ( !check_ref( $slice, 'Bio::EnsEMBL::Slice' )
            && !check_ref( $slice, 'Bio::EnsEMBL::LRGSlice' ) )
        {
            throw('slice argument must be a Bio::EnsEMBL::Slice');
        }
        $self->{ 'slice' } = $slice;
    } elsif ( @_ > 1 ) {
        delete($self->{ 'slice' });
    }
    return $self->{ 'slice' };
}
```

Reimplemented in [Bio::EnsEMBL::Exon](#), and [Bio::EnsEMBL::Map::DitagFeature](#).

# Ensembl Core DB Documentation (1)

Go to [http://www.ensembl.org/info/docs/api/core/core\\_schema.html](http://www.ensembl.org/info/docs/api/core/core_schema.html)

## Ensembl Core - Schema documentation

This document gives a high-level description of the tables that make up the EnsEMBL core schema. Tables are grouped into logical groups, and the purpose of each table is explained. It is intended to allow people to familiarise themselves with the schema when encountering it for the first time, or when they need to use some tables that they've not used before.

This document refers to version **73** of the EnsEMBL core schema.

The core database schema is available in several diagrams (PDF format) here:



## List of the tables:

# Ensembl Core DB Documentation (2)

Scroll down to the list of Fundamental Tables and click on the gene table.

■ **Fundamental Tables**

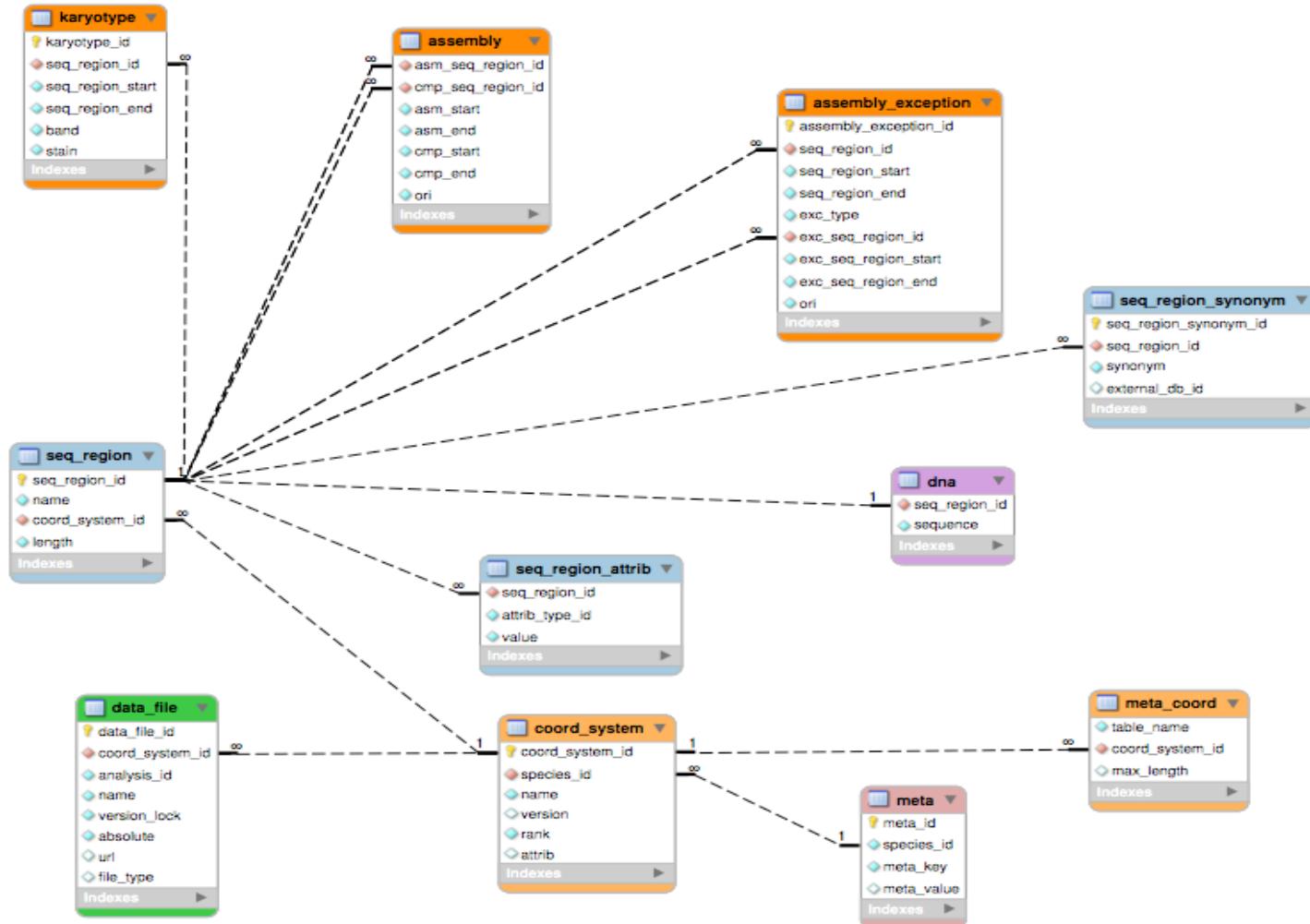
- [alt allele](#)
- [alt allele attrib](#)
- [alt allele group](#)
- [analysis](#)
- [analysis description](#)
- [attrib type](#)
- [dna\\_align\\_feature](#)
- [exon](#)
- [exon\\_transcript](#)
- [gene](#)
- [gene\\_attrib](#)
- [protein\\_align\\_feature](#)
- [protein\\_feature](#)
- [splicing\\_event](#)
- [splicing\\_event\\_feature](#)
- [splice transcript pair](#)
- [supporting feature](#)
- [transcript](#)
- [transcript\\_attrib](#)
- [transcript supporting feature](#)
- [translation](#)
- [translation\\_attrib](#)

# Ensembl Core DB Documentation (3)

Click on the ‘show columns’ link to the right to expand a list of gene table columns, their types, descriptions and indices over the columns.

gene		<a href="#">Hide columns</a>   <a href="#">[Back to top]</a>	
Column	Type	Default value	Description
gene_id	INT(10)		Primary key, internal identifier.
biotype	VARCHAR(40)		Biotype, e.g. protein_coding.
analysis_id	SMALLINT		Foreign key references to the <a href="#">analysis</a> table.
seq_region_id	INT(10)		Foreign key references to the <a href="#">seq_region</a> table.
seq_region_start	INT(10)		Sequence start position.
seq_region_end	INT(10)		Sequence end position.
seq_region_strand	TINYINT(2)		Sequence region strand: 1 - forward; -1 - reverse.
display_xref_id	INT(10)		External reference for Ensembl web site. Foreign key re e.g ensembl, havana etc.
source	VARCHAR(20)		
status	ENUM('KNOWN', 'NOVEL', 'PUTATIVE', 'PREDICTED', 'KNOWN_BY_PROJECTION', 'UNKNOWN', 'ANNOTATED')		Status, e.g.'KNOWN', 'NOVEL', 'PUTATIVE', 'PREDICTED' 'KNOWN_BY_PROJECTION', 'UNKNOWN'.
description	TEXT		Gene description

# Ensembl Core DB Documentation (4)



# Exercise 1

- a) Find documentation for the Exon class in the Ensembl core code base. Which method would you use to retrieve the DNA sequence for an exon? What is the return type for this method?
  
- b) Can you find a table which stores stable ids for transcripts? Which table stores DNA sequence? How many columns does this table have?

# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. **The Registry & Ensembl API script design**
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references

# The Registry

We know how to use Data Objects and Object Adaptors.

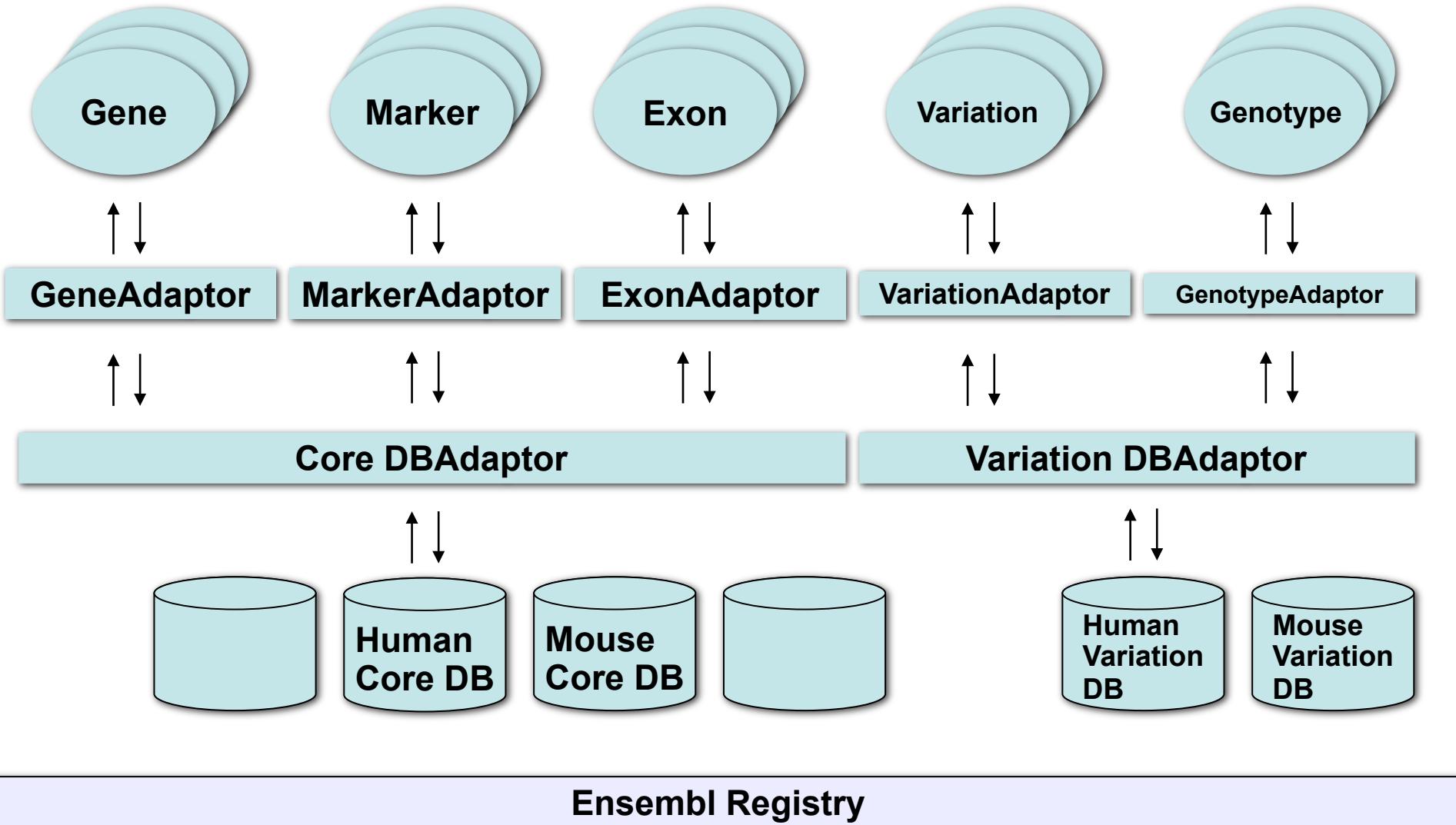
How do we make sure we get those from the right database?

This is what we use the Registry for.

The Registry:

- loads all databases of the same version as the API
- lazy loads so no connections are made until requested

# Ensembl API Architecture



# The beginning of each script

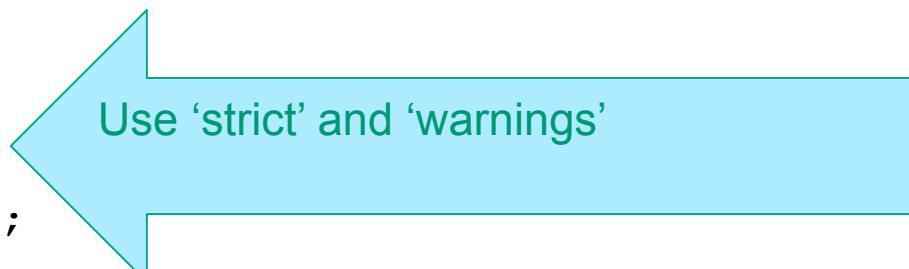
```
#!/usr/bin/env perl

use strict;
use warnings;
use Bio::EnsEMBL::Registry;

my $registry = 'Bio::EnsEMBL::Registry';

$registry->load_registry_from_db(
    -host => 'ensembldb.ensembl.org',
    -user => 'anonymous'
);

my $gene_adaptor = $registry->get_adaptor('Human', 'Core', 'Gene');
```



↑  
species    database    object type

# Connecting to Ensembl Genomes

```
#!/usr/bin/env perl

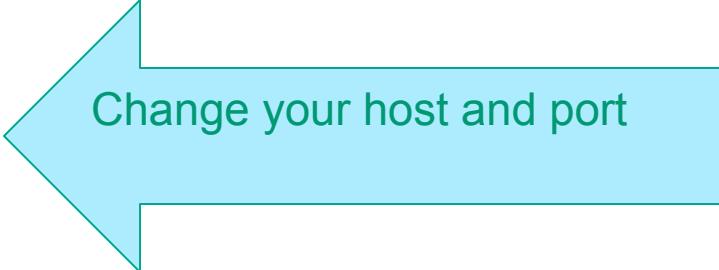
use strict;
use warnings;
use Bio::EnsEMBL::Registry;

my $registry = 'Bio::EnsEMBL::Registry';

$registry->load_registry_from_db(
    -host  => 'mysql.ebi.ac.uk',
    -port  => 4157,
    -user  => 'anonymous'
);

my $gene_adaptor = $registry->get_adaptor('silkmoth', 'Core', 'Gene');
```

↑  
species    database    object type



Change your host and port

# Exercise 2

Create a script which uses the method `load_registry_from_db` to load all databases into the Registry and prints the names of the databases loaded.

*Hint: Have a look at the Doxygen documentation for the Registry object and method `load_registry_from_db` (<http://www.ensembl.org/info/docs/Doxygen/index.html>).*

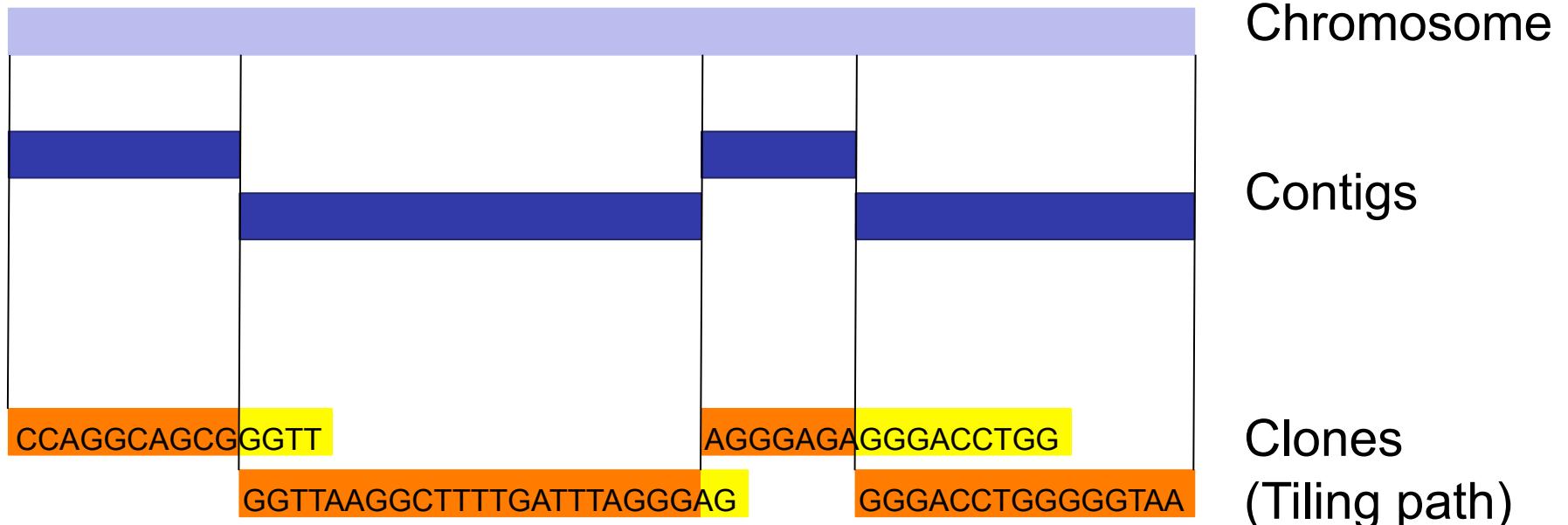
# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. **Coordinate systems & slices**
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references

# Coordinate Systems (1)

Ensembl stores features and DNA sequence in a number of coordinate systems.

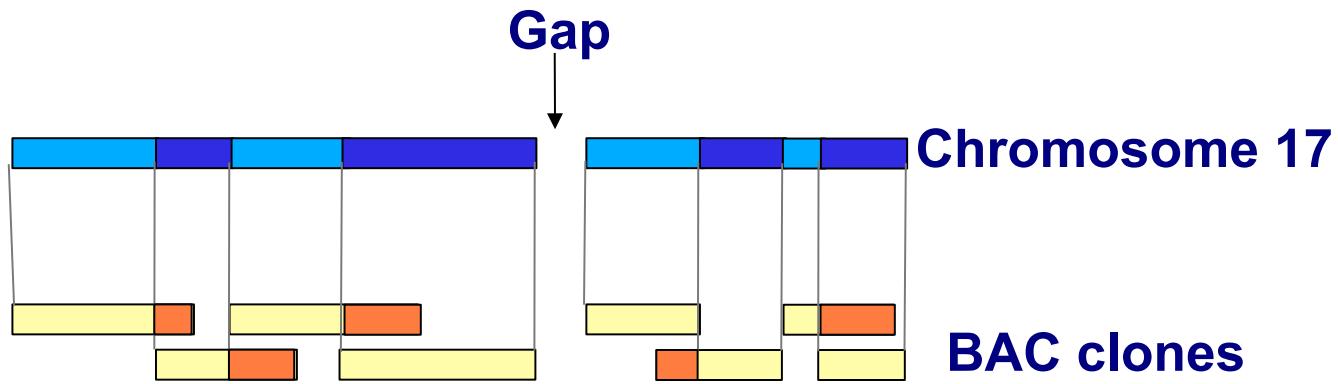
Top level coordinate system



Sequence level coordinate system

# Coordinate Systems (2)

Regions in one coordinate system may be constructed from a tiling path of regions from another coordinate system.



# Coordinate Systems - Code Example

```
# Obtain all coordinate systems for human

use Bio::EnsEMBL::Registry;
my $registry = 'Bio::EnsEMBL::Registry';

$registry->load_registry_from_db(
    -host => 'ensembldb.ensembl.org',
    -user => 'anonymous'
);
my $coordsystem_adaptor = $reg->get_adaptor( 'Human', 'Core', 'CoordSystem' );

my $coordsystems = $coordsystem_adaptor->fetch_all;

while ( my $coordsystem = shift @{$coordsystems} ) {
    print $coordsystem->name, "\t",
        $coordsystem->version, "\t",
        $coordsystem->rank ,"\n";
}
```



Note use of 'while' and 'shift' instead of 'foreach' – more memory efficient way for large datasets

# Coordinate Systems – Code Output

## OUTPUT:

name	version	rank
chromosome	GRCh37	1
supercontig	GRCh37	2
clone		3
contig		4
chromosome	NCBI36	5
chromosome	NCBI35	6
chromosome	NCBI34	7
lrg		8

Latest assembly, top level

Latest assembly, sequence level

Old assemblies, used for mapping features between assembly versions.

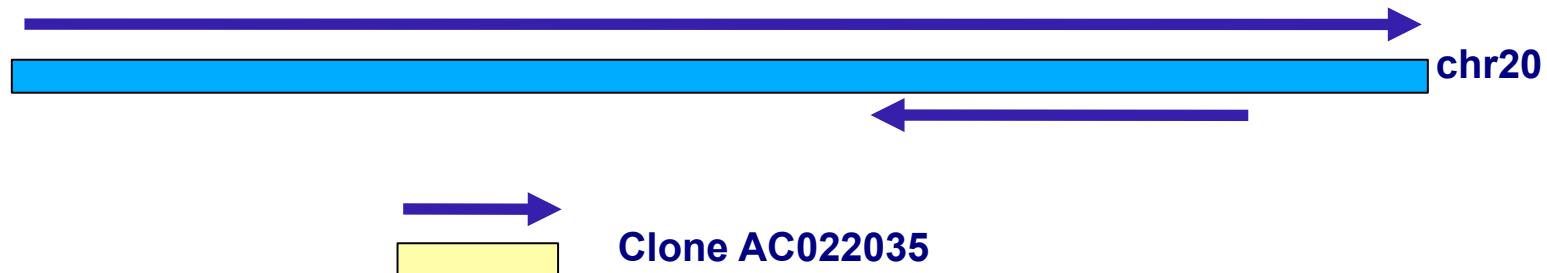
Locus-Reference Genes, used in clinical tests

# Slices

A *Slice Data Object* represents an arbitrary region of a genome, a slice of a Sequence Region.

*Slices* are not directly stored in the database.

A *Slice* is used to request sequence or features from a specific region in a specific coordinate system.



# Slices - Code Example (1)

```
# get a slice covering the entire human Y chromosome

my $slice_adaptor = $registry->get_adaptor( 'Human', 'Core', 'Slice' );
my $slice = $slice_adaptor->fetch_by_region( 'chromosome', 'Y' );
print "Coord system:\t", $slice->coord_system_name, "\n",
      "Seq region:\t", $slice->seq_region_name, "\n",
      "Start:\t\t", $slice->start, "\n",
      "End:\t\t", $slice->end, "\n",
      "Strand:\t\t", $slice->strand, "\n",
      "Slice:\t\t", $slice->name, "\n";
```

## OUTPUT:

Coord system: chromosome  
Seq region: Y  
Start: 1  
End: 59373566  
Strand: 1  
Slice: chromosome:GRCh37:Y:1:59373566:1

# Slices - Code Example (2)

```
# get the slice adaptor
$slice_adaptor = $reg->get_adaptor('human', 'core', 'slice');

# fetch a slice on a region of chromosome 12
$slice = $slice_adaptor->fetch_by_region('chromosome', '12',
                                             1e6, 2e6);

# print out the sequence from this region
print $slice->seq, "\n";

# get all clones in the database and print out their names
@slices = @{} $slice_adaptor->fetch_all('clone') {};
foreach $slice (@slices) {
    print $slice->seq_region_name, "\n";
}
```

# Exercise 3

- (a) Fetch all chromosomes for human. Determine their number and print the name and length for each of them. *The number of chromosomes is probably not what you would expect! Why is this?*
- (b) Use the gene stable id 'ENSG00000101266' to fetch a slice surrounding this gene with 2kb of flanking sequence.  
(a) hint: use the Ensembl API documentation to find an appropriate method in SliceAdaptor class which retrieves a slice given a gene stable id (pay attention to the method's arguments)
- (c) Fetch the sequence of the first 10MB of chromosome 20 and write it to a file in FASTA format. Print the number of genes in this region.  
✓ hint: Slice objects inherit from Bio::Seq so can be written to file easily using Bio::SeqIO, e.g.:

```
my $output = Bio::SeqIO->new( -file=>'>filename.fasta',
    -format=>'Fasta');
$output->write_seq($slice);
```

# Outline

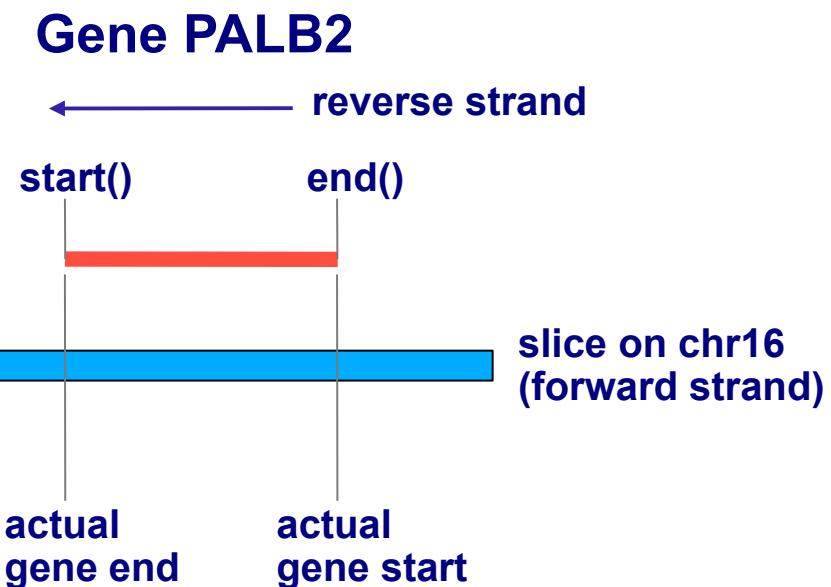
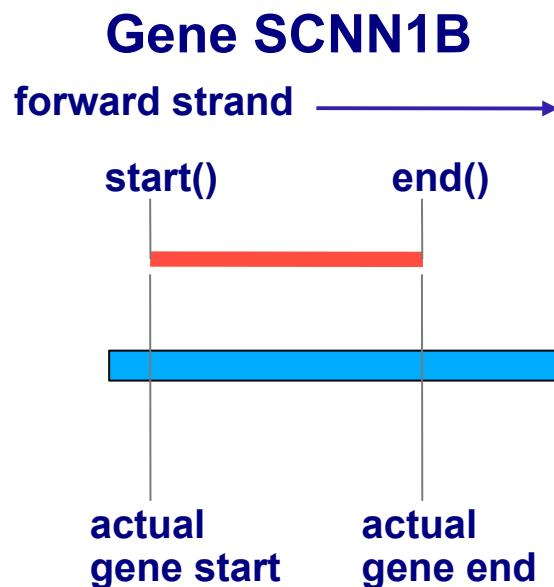
- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. **Features**
- g. Genes, transcripts, exons & translations
- h. External references

# Features (1)

Features have a defined location on the genome and are stored in a single coordinate system

All Features have a *start*, *end*, *strand* and *slice*

*Start* and *end* are plotted onto the forward strand:  $\text{start} < \text{end}$

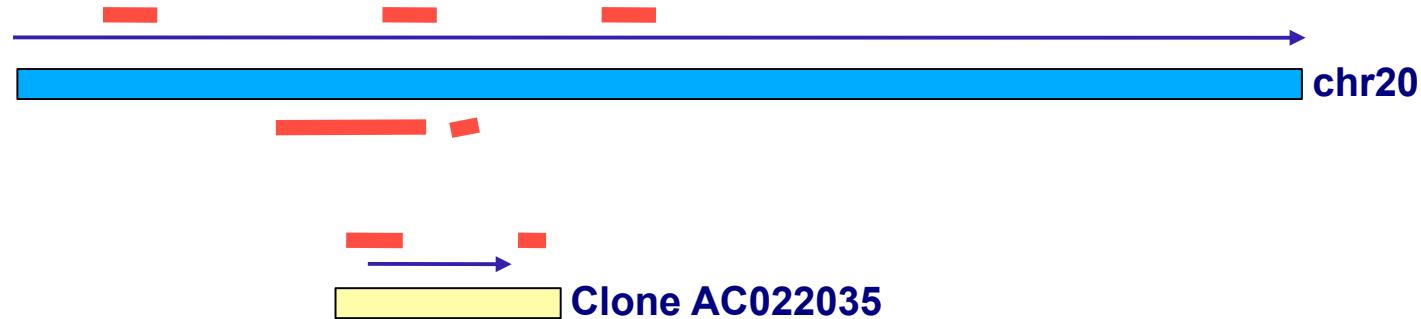


# Features (2)

*slice* method returns the Slice object with which the *Feature* is associated

*feature\_Slice* method returns the *Slice* object which covers the *Feature*

*Features* are retrieved from Object Adaptors using identifiers or regions (slices).



# Features Objects – Biological Correspondence

Object	Biological entity
<b>Gene, Transcript, Exon</b>	<b>Ensembl gene models</b>
PredictionTranscript, PredictionExon	Genscan gene models
DNAAlignFeature, ProteinAlignFeature	cDNAs, proteins
<b>RepeatFeature</b>	<b>repeats</b>
MarkerFeature	markers
OligoFeature	microarray probes
KaryotypeBandFeature	cytogenetic bands
SimpleFeature	results of cpg, Eponine, FirstEF and tRNAscan
MiscFeature	clones, ENCODE regions
ProteinFeature	protein domains

# Exercise 4

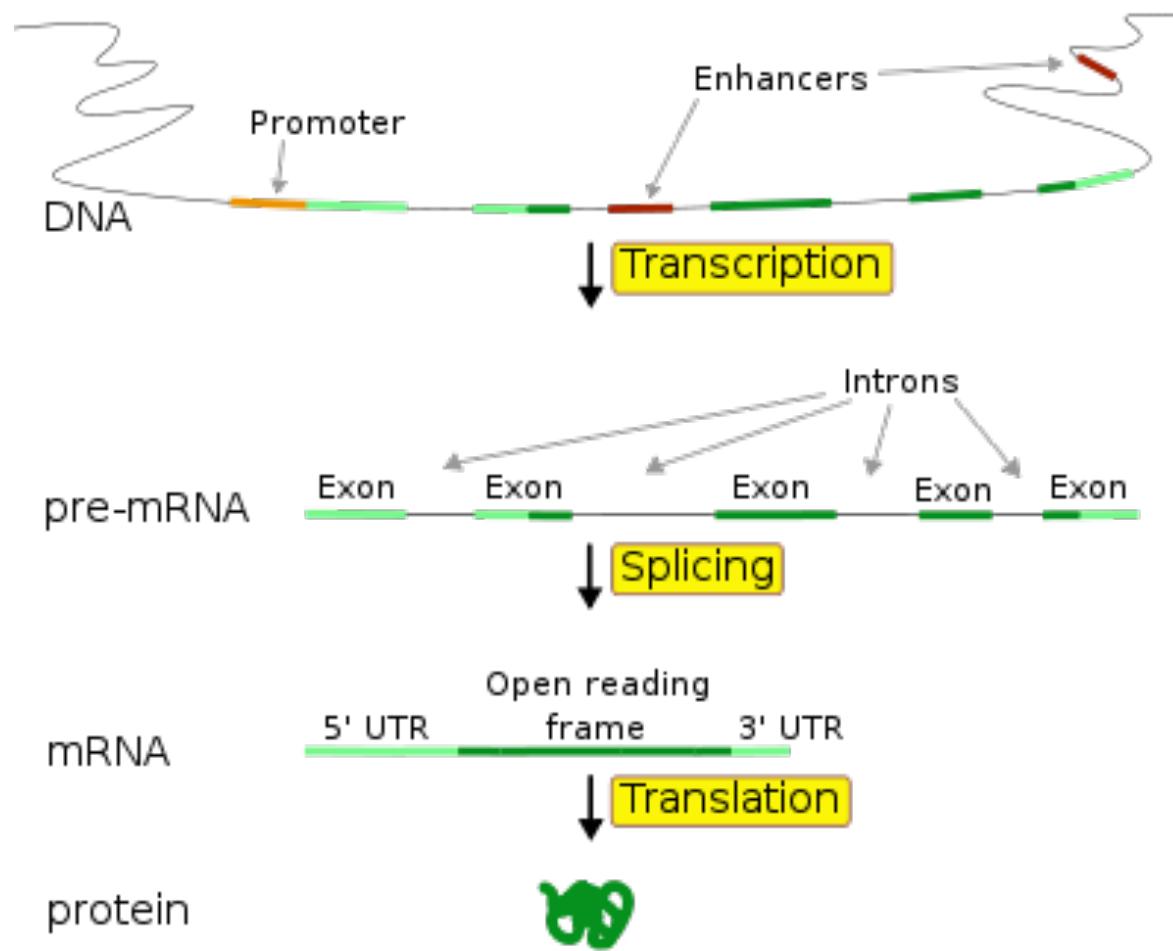
- (a) Get all the repeat features from chromosome 20:1-500k. Print out the name and position of each on the chromosome and the total number.
  - hint: create a slice, retrieve repeat features on this slice
- (b) Find which genomic region the RefSeq dna entry NM\_000059.3 was mapped to. Print the name of the region and coordinates of the alignment on the genome as well as the name of the region and coordinates of the alignment on the RefSeq dna entry. Print the score and percentage identity for the alignment.
  - ✓ hint: use DnaAlignFeatureAdaptor; use the core schema documentation as a guide to appropriate methods which correspond to columns in dna\_align\_feature table

*A list of useful Feature methods is in the Appendix at the end of the presentation slides*

# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations**
- h. External references

# Genes, transcripts, exons & translations



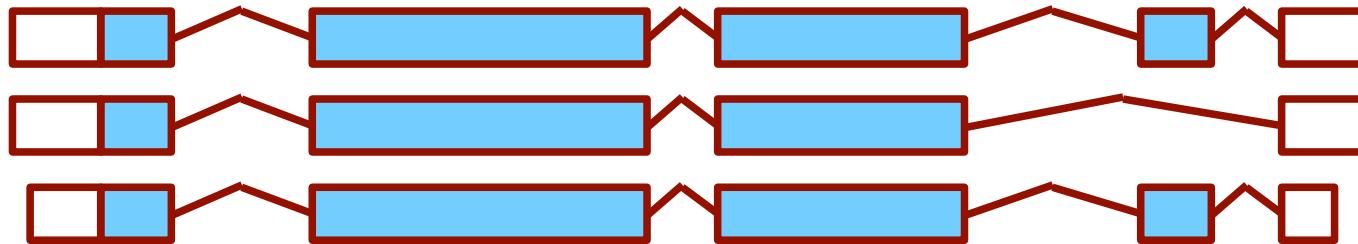
# Genes, Transcripts and Exons

Genes, Transcript and Exons are objects that can be used just like any other Feature object

A Gene is a set of alternatively spliced Transcripts

A Transcript is a set of Exons

Introns are not explicitly defined in the database



# Genes, Transcripts, Exons – Code Output

```
# helper function: returns location and stable_id string for a feature
sub get_string {
    my $feature = shift;
    my $stable_id = $feature->stable_id;
    my $seq_region = $feature->slice->seq_region_name;
    my $start = $feature->start;
    my $end = $feature->end;
    my $strand = $feature->strand;
    return "$stable_id $seq_region:$start-$end($strand)";
}

# fetch a gene by its stable identifier
my $gene = $gene_adaptor->fetch_by_stable_id('ENSG00000123427');

# print out the gene, its transcripts, and its exons
print "Gene: ", get_string($gene), "\n";
foreach my $transcript (@{ $gene->get_all_Transcripts }) {
    print " Transcript: ", get_string($transcript), "\n";
    foreach my $exon (@{ $transcript->get_all_Exons }) {
        print "     Exon: ", get_string($exon), "\n";
    }
}
```

# Genes, Transcripts, Exons – Code Output

## OUTPUT:

Gene: ENSG00000123427 12:58165275-58176324(1)

Transcript: ENST00000548256 12:58165275-58174213(1)

Exon: ENSE00002360002 12:58165275-58165408(1)

Exon: ENSE00003090890 12:58166800-58166911(1)

Exon: ENSE00003030714 12:58168412-58168550(1)

Exon: ENSE00002406112 12:58174038-58174213(1)

Transcript: ENST00000551420 12:58165983-58174563(1)

Exon: ENSE00002355737 12:58165983-58166180(1)

Exon: ENSE00003223131 12:58166800-58166911(1)

Exon: ENSE00002376752 12:58174038-58174563(1)

Transcript: ENST00000300209 12:58166383-58176186(1)

Exon: ENSE00002301479 12:58166383-58166684(1)

Exon: ENSE00003090890 12:58166800-58166911(1)

Exon: ENSE00002393444 12:58174038-58176186(1)

etc.

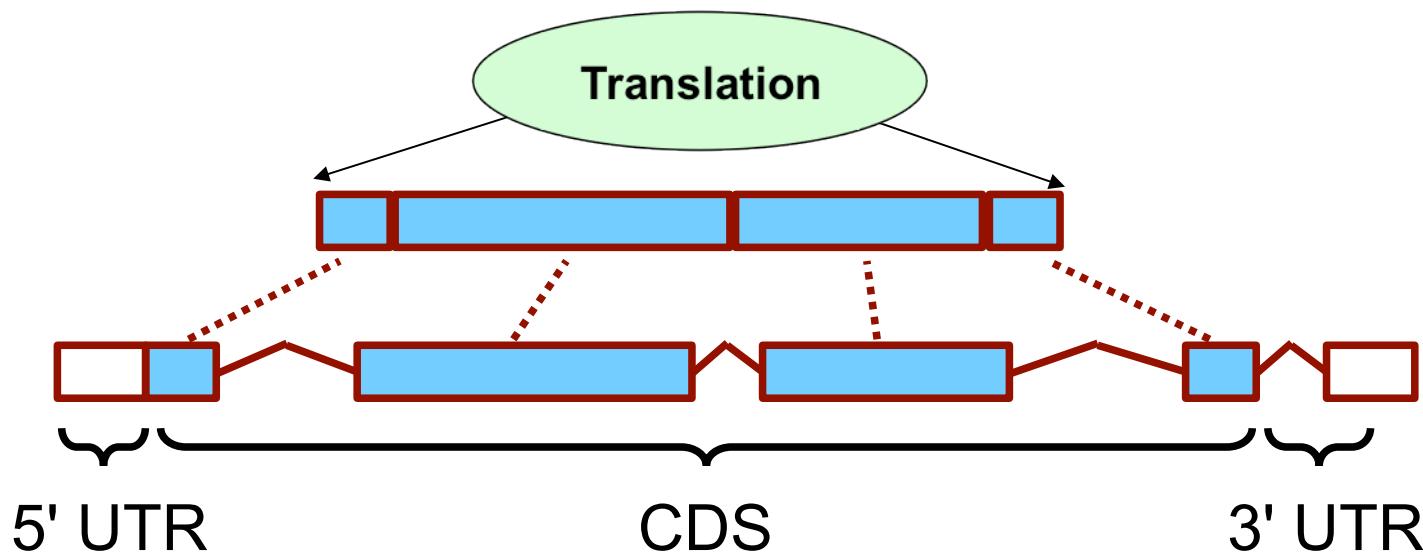
# Translations

Translations are not Features.

A Translation object defines the UTR and CDS of a Transcript.

Peptides are not stored in the database, they are computed on the fly using Transcript objects.

Not all transcripts have a translation (e.g. ncRNAs)



# Translations – Code Example

```
my $transcript_adaptor = $registry->get_adaptor( 'Homo sapiens', 'Core',
    'Transcript' );
# fetch a transcript from the database
my $transcript =
    $transcript_adaptor->fetch_by_stable_id('ENST00000333012');

# obtain the translation of the transcript
my $translation = $transcript->translation;

# print out the translation info
print "Translation: ", $translation->stable_id, "\n";
print "Start Exon: ", $translation->start_Exon->stable_id, "\n";
print "End Exon:    ", $translation->end_Exon->stable_id, "\n";

# cDNA start and end (spliced sequence with UTR)
print "Start : ", $translation->cdna_start, "\n";
print "End   : ", $translation->cdna_end, "\n";

# print the peptide which is the product of the translation
print "Peptide : ", $transcript->translate->seq, "\n";
```

# Translations – Code Output

## OUTPUT:

Translation: ENSP00000327425

Start Exon: ENSE00002340145

End Exon: ENSE00002428887

Start : 48

End : 497

Peptide :

MADPGPDPESESESVFPREVGLFADSYSEKSQFCFCGHVLTITQNFGSRLGVAARVWDAALSLCNYFESQNVDFR  
GKKVIELGAGTGIVGILAALQGAYGLVRETEDDVICEQELWRGMRGACGHALSMSTMTPWESIKGSSVRGGCYHH

# Exercise 5

(a) Fetch gene 'CSNK2A1' and print the number of its transcripts and exons.

- ✓ hint: use GeneAdaptor method `fetch_by_display_label`; remember that not all transcripts have a translation

Gene: CSNK2A1 ENSG00000101266

Description casein kinase 2, alpha 1 polypeptide [Source:HGNC Symbol;Acc:2457]  
Location Chromosome 20: 459,116-524,482 reverse strand.  
Transcripts This gene has 5 transcripts

Show/hide columns Filter

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
CSNK2A1-001	<a href="#">ENST00000349736</a>	4299	<a href="#">ENSP00000339247</a>	391	Protein coding	<a href="#">CCDS13003</a>
CSNK2A1-002	<a href="#">ENST00000217244</a>	4416	<a href="#">ENSP00000217244</a>	391	Protein coding	<a href="#">CCDS13003</a>
CSNK2A1-201	<a href="#">ENST00000400217</a>	2507	<a href="#">ENSP00000383076</a>	255	Protein coding	<a href="#">CCDS13004</a>
CSNK2A1-202	<a href="#">ENST00000400227</a>	1483	<a href="#">ENSP00000383086</a>	385	Protein coding	-
CSNK2A1-003	<a href="#">ENST00000460062</a>	768	No protein product	-	Processed transcript	-

(b) For the above gene, get all the transcripts and list the number of exons in each and the translations.

(c) Why do the exon numbers not match?

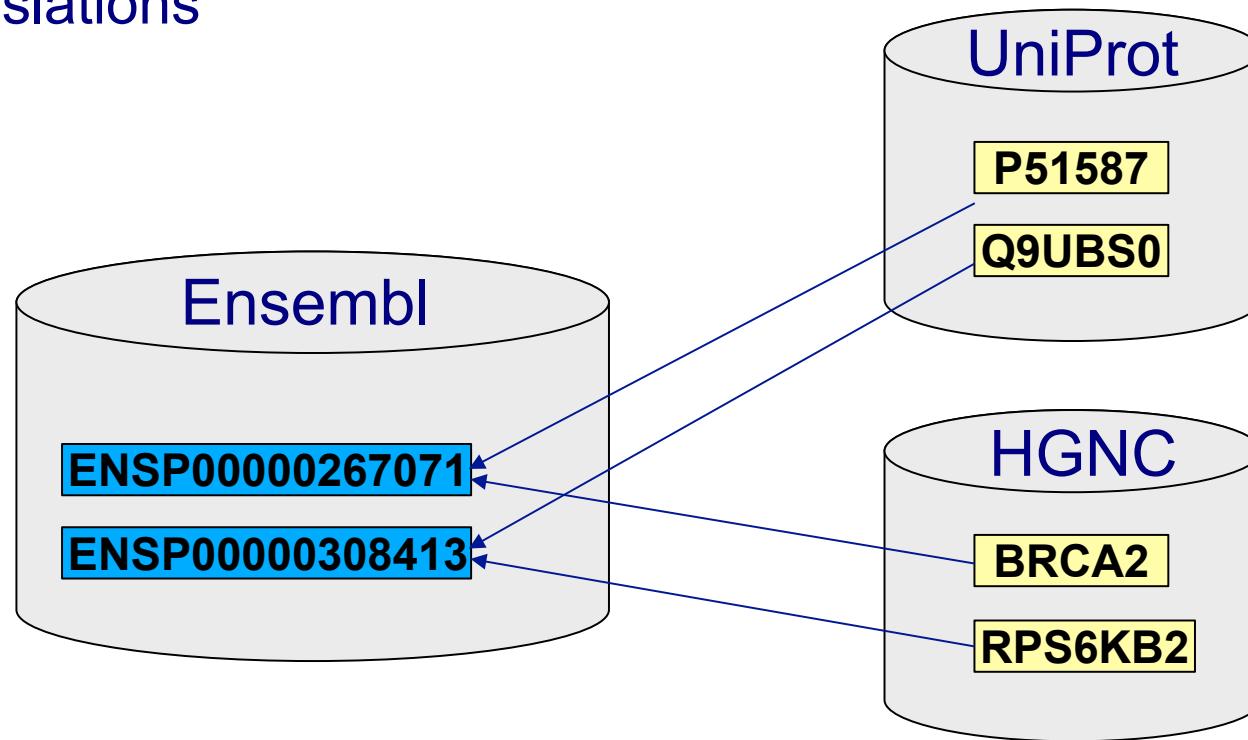
# Outline

- a. Introduction
- b. Data objects & object adaptors
- c. Ensembl documentation
- d. The Registry & Ensembl API script design
- e. Coordinate systems & slices
- f. Features
- g. Genes, transcripts, exons & translations
- h. External references**

# External References

Ensembl cross references its Gene models with identifiers from other databases, such as HGNC, WikiGenes, UniProtKB/Swiss-Prot, RefSeq, MIM etc.

External References (Xrefs) can be linked to genes, transcripts or translations



# Xrefs - Code Example (1)

```
# Obtain external references for Ensembl gene ENSG00000139618
```

```
my $gene = $gene_adaptor->fetch_by_stable_id( 'ENSG00000139618' );
```

```
my $gene_xrefs = $gene->get_all_DBEntries;  
print "Xrefs on gene level: \n\n";
```

```
foreach my $gene_xref( @{$gene_xrefs} ) {  
    print $gene_xref->dbname, ":", $gene_xref->display_id, "\n";  
}
```

```
my $all_xrefs = $gene->get_all_DBLinks;
```

```
print "\nXrefs on gene, transcript and protein level: \n\n";  
foreach my $all_xref( @{$all_xrefs} ) {  
    print $all_xref->dbname, ":", $all_xref->display_id, "\n";  
}
```

this method will only return xrefs linked to the object it's called on (e.g. gene)

this method will return xrefs on all levels (gene, transcript and translation)

# Xrefs – Code Output (1)

## OUTPUT:

Xrefs on gene level:

OTTG:OTTHUMG00000017411  
ArrayExpress:ENSG00000139618  
EntrezGene:BRCA2  
HGNC:BRCA2  
MIM\_MORBID: BREAST CANCER [#114480]  
MIM\_MORBID: FANCONI ANEMIA, COMPLEMENTATION GR [#227650]  
MIM\_MORBID: FANCONI ANEMIA, COMPLEMENTATION GR [#605724]  
MIM\_MORBID: BREAST-OVARIAN CANCER, FAMILIAL, S [#612555]  
MIM\_MORBID: GLIOMA SUSCEPTIBILITY 3 [#613029]  
MIM\_MORBID: PANCREATIC CANCER, SUSCEPTIBILITY [#613347]  
MIM\_GENE: BRCA2 GENE [\*]  
Orphanet:Hereditary breast and ovarian cancer syndrome  
Orphanet:Familial prostate cancer  
Orphanet:Fanconi anemia  
Orphanet:Familial pancreatic carcinoma  
UniGene:Hs.34012  
Uniprot\_genename:BRCA2  
WikiGene:BRCA2

DBASS3:BRCA2

DBASS5:BRCA2

Xrefs on gene, transcript and translation level:

(same as on gene level + transcript and translation level)

CCDS:CCDS9344.1

Uniprot/SWISSPROT:BRCA2\_HUMAN

RefSeq\_peptide:NP\_000050.2

RefSeq\_mRNA:NM\_000059.3

Vega\_transcript:BRCA2-001

Vega\_transcript:OTTHUMT00000046000

Uniprot/SPTREMBL:E9PIQ1\_HUMAN

Uniprot/SPTREMBL:K4JTT2\_HUMAN

Uniprot/SPTREMBL:K4JXT6\_HUMAN

Uniprot/SPTREMBL:K4K7W0\_HUMAN

Uniprot/SPTREMBL:Q8IU64\_HUMAN

Uniprot/SPTREMBL:Q8IU77\_HUMAN

Uniprot/SPTREMBL:Q8IU82\_HUMAN

Uniprot/SPTREMBL:Q9H4L3\_HUMAN

UCSC:uc001uub.

etc.

# Xrefs - Code Example (2)

```
# Retrieve Ensembl IDs for a list of UniProt protein IDs

my $translation_adaptor
  = $registry->get_adaptor(
    'Human', 'Core', 'Translation');

my @uniprot_ids = qw(P51587 P15056 B8A597 B8A595 B7ZW72);

foreach my $uniprot_id (@uniprot_ids){

  my @trans = @{
    $translation_adaptor->fetch_all_by_external_name($uniprot_id, 'Uniprot%')
  };

  foreach my $translation (@trans){
    print $translation->stable_id."\t".$uniprot_id."\n";
  }
}
```

Proteins map to Ensembl Translation objects so we will use a TranslationAdaptor

# Xrefs – Code Output (2)

## OUTPUT :

ENSP00000369497	P51587
ENSP00000439902	P51587
ENSP00000288602	P15056
ENSP00000387217	B8A597
ENSP00000386781	B8A595
ENSP00000442513	B7ZW72
ENSP00000307640	B7ZW72

Cross references can map to more than one Ensembl identifier

# Exercise 6

Retrieve a list of GO term IDs and term names linked to the gene with stable id ‘ENSG00000139618’

- ✓ Use `get_all_DBLinks` with an external database name argument to restrict the number of xrefs returned
- ✓ Ontology term data such as term name and definition are stored outside of the core database. Create an `OntologyTerm` Adaptor with the help of the Registry method `get_adaptor` using arguments: ‘Multi’ (species), ‘Ontology’ (database type), ‘`OntologyTerm`’ (adaptor type)
- ✓ For all xrefs returned by `get_all_DBLinks` use the `OntologyTerm` Adaptor to fetch the relevant term and print its accession and name (xref `display_id` is the same as term accession)

# Recap - Ensembl API script design

Always:

- Load the registry

Which features (genes, repeats, SNPs, etc.) are in my particular region of interest?

- Get the SliceAdaptor
- Fetch the Slice for your region of interest
- Get the features from your Slice

What do we know about a particular gene (or any other feature)?

- Get the GeneAdaptor
- Fetch your Gene of interest
- Get more details about the gene:
  - Gene structure (transcripts, exons, translations)
  - Annotations: GO xrefs, HGNC symbols, etc.
  - Features in the same region -> get the Slice for the Gene!

# Documentation & Help

- Installation instructions, web-browsable version of the POD (Perldoc), database schema and tutorial:  
<http://www.ensembl.org/info/docs/api/index.html>
- Inline Perl POD (Plain Old Documentation)
- [dev@ensembl.org](mailto:dev@ensembl.org) mailing list:  
<http://www.ensembl.org/info/about/contact/mailing.html>  
searchable mailing list archive:  
<http://blog.gmane.org/gmane.science.biology.ensembl.devel>
- Ensembl helpdesk:  
[helpdesk@ensembl.org](mailto:helpdesk@ensembl.org)

# Citing Ensembl

D48–D55 *Nucleic Acids Research*, 2013, Vol. 41, Database issue  
doi:10.1093/nar/gks1236

Published online 30 November 2012

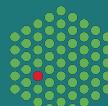
## Ensembl 2013

**Paul Flicek<sup>1,2,\*</sup>, Ikhlaq Ahmed<sup>1</sup>, M. Ridwan Amode<sup>2</sup>, Daniel Barrell<sup>2</sup>, Kathryn Beal<sup>1</sup>, Simon Brent<sup>2</sup>, Denise Carvalho-Silva<sup>1</sup>, Peter Clapham<sup>2</sup>, Guy Coates<sup>2</sup>, Susan Fairley<sup>2</sup>, Stephen Fitzgerald<sup>1</sup>, Laurent Gil<sup>1</sup>, Carlos García-Girón<sup>2</sup>, Leo Gordon<sup>1</sup>, Thibaut Hourlier<sup>2</sup>, Sarah Hunt<sup>1</sup>, Thomas Juettemann<sup>1</sup>, Andreas K. Kähäri<sup>2</sup>, Stephen Keenan<sup>1</sup>, Monika Komorowska<sup>1</sup>, Eugene Kulesha<sup>1</sup>, Ian Longden<sup>1</sup>, Thomas Maurel<sup>1</sup>, William M. McLaren<sup>1</sup>, Matthieu Muffato<sup>1</sup>, Rishi Nag<sup>2</sup>, Bert Overduin<sup>1</sup>, Miguel Pignatelli<sup>1</sup>, Bethan Pritchard<sup>2</sup>, Emily Pritchard<sup>1</sup>, Harpreet Singh Riat<sup>2</sup>, Graham R. S. Ritchie<sup>1</sup>, Magali Ruffier<sup>1</sup>, Michael Schuster<sup>1</sup>, Daniel Sheppard<sup>2</sup>, Daniel Sobral<sup>1</sup>, Kieron Taylor<sup>1</sup>, Anja Thormann<sup>1</sup>, Stephen Trevanion<sup>2</sup>, Simon White<sup>2</sup>, Steven P. Wilder<sup>1</sup>, Bronwen L. Aken<sup>2</sup>, Ewan Birney<sup>1</sup>, Fiona Cunningham<sup>1</sup>, Ian Dunham<sup>1</sup>, Jennifer Harrow<sup>2</sup>, Javier Herrero<sup>1</sup>, Tim J. P. Hubbard<sup>2</sup>, Nathan Johnson<sup>1</sup>, Rhoda Kinsella<sup>1</sup>, Anne Parker<sup>2</sup>, Giulietta Spudich<sup>1</sup>, Andy Yates<sup>1</sup>, Amonida Zadissa<sup>2</sup> and Stephen M. J. Searle<sup>2</sup>**

<sup>1</sup>European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton Cambridge CB10 1SD, UK and

<sup>2</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

Received October 11, 2012; Revised October 31, 2012; Accepted November 1, 2012



# The Ensembl Team



# Ensembl Acknowledgements

## Ensembl Team & Others

### Ensembl Core

Andy Yates

Stephen Keenan

Kieron Taylor

Alessandro Vullo

Paul Flicek, Steve Searle and the  
entire Ensembl Team

## Funding



European Commission  
Framework Programme 7



# Appendix – CoordSystem Methods

Attribute	Example value(s)	Method(s)
name	chromosome, scaffold, contig, clone	\$coordsystem->name
version	GRCh37, NCBI36, NCBIM37	\$coordsystem->version

# Appendix – Feature Methods

Attribute	Example value(s)	Method(s)
name	AluSp, D1S2217	\$feature->display_id
coordinates		\$feature->seq_region_name \$feature->start \$feature->end \$feature->seq_region_start \$feature->seq_region_end \$feature->strand
sequence		\$feature->seq
length	399	\$feature->length
slice	returns Slice object with which feature is associated	\$feature->slice
feature slice	returns Slice object that covers feature	\$feature->feature_Slice

# Appendix – Gene Methods

Attribute	Example value(s)	Method(s)
stable ID	ENSG00000139618	\$gene->stable_id
name	BRCA2	\$gene->external_name
description	breast cancer 2, early onset	\$gene->description
biotype	protein_coding, miRNA	\$gene->biotype
analysis	ensembl, havana, ensembl_havana_gene	\$gene->analysis->logic_name
status	KNOWN, NOVEL	\$gene->status
transcripts	returns listref of Transcript objects	\$gene->get_all_Transcripts
exons	returns listref of Exon objects	\$gene->get_all_Exons

# Appendix - Transcript Methods

Attribute	Example value(s)	Method(s)
stable ID	ENST00000380152	\$transcript->stable_id
name	BRCA2-001	\$transcript->external_name
biotype	protein_coding, nonsense-mediated_decay	\$transcript->biotype
analysis	ensembl, havana, ensembl_havana_transcript	\$transcript->analysis-> logic_name
status	KNOWN, NOVEL	\$transcript->status
CDS (spliced sequence, no UTR)	ATGCCTATTGGATCCAAAGAGAGGC...	\$transcript->translateable_seq
UTRs	returns Seq object	\$transcript->five_prime_utr \$transcript->three_prime_utr
cDNA (spliced sequence + UTR)	GGGCTTGTGGCGCGAGCTTCTGAAA...	\$transcript->spliced_seq
translation	returns Translation object	\$transcript->translation
exons	returns listref of Exon objects	\$transcript->get_all_Exons
introns	returns listref of Intron objects	\$transcript->get_all_Introns

# Appendix - Translation Methods

Attribute	Example value(s)	Method(s)
stable id	ENSP0000369497	\$translation->stable_id
length	3418	\$translation->length
sequence	MPIGSKERPTFFEIFKTRCNKADLG...	\$translation->seq