



Ensembl Variation API

Hinxton, 26th of January 2016

Laurent Gil
Ensembl Variation

lgil@ebi.ac.uk

Materials:
<http://www.ebi.ac.uk/~lgil/workshops/>

Schedule

09:30 - 10:45 Variation

10:45 - 11:15 Coffee break

11:15 - 12:30 Variation

12:30 - 13:30 Lunch break

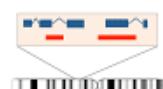
13:30 - 15:30 Variation

15:30 - 16:00 Tea break

16:00 - 17:00 Variation + VEP

Important URLs

- Links for the course material
 - http://www.ebi.ac.uk/~lgil/workshops/2016-01_Hinxton/
 - Ensembl **Variation** API documentation
 - <http://www.ensembl.org/info/docs/Doxygen/variation-api/index.html>
 - Ensembl **Core** API documentation
 - <http://www.ensembl.org/info/docs/Doxygen/core-api/index.html>
 - Variation data documentation:
 - <http://www.ensembl.org/info/genome/variation/index.html>



Genomic context



Genes and regulation



Population genetics



Sample genotypes



Linkage disequilibrium



Phenotype data



Citation



Phylogenetic context

ATTCATT
CGG**S**GTG
TCATGCT

Flanking sequence



Overview

Introduction to the Variation data and API

Key data/objects

- Variation
- Allele
- Genotypes
- Variation feature (location of the variants)
- Structural variation and structural variation feature (location)

Additional information/annotation

- Variation set
- Phenotype data
- Publication data

Analysis

- Variation consequences
- Linkage disequilibrium

Variant Effect Predictor (VEP)

An introduction to genomic variation

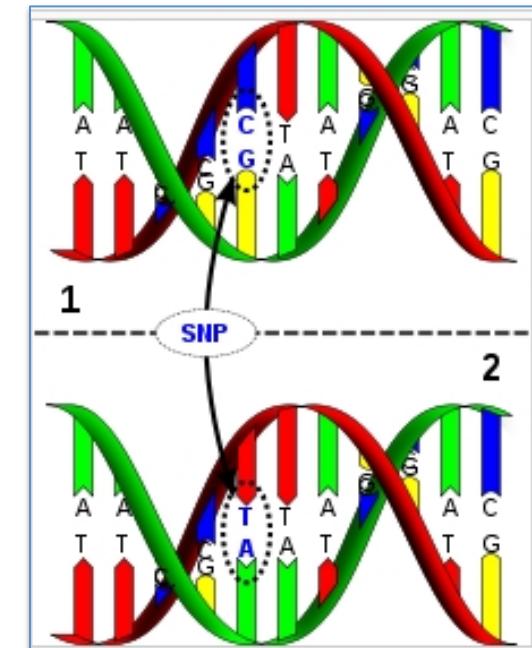
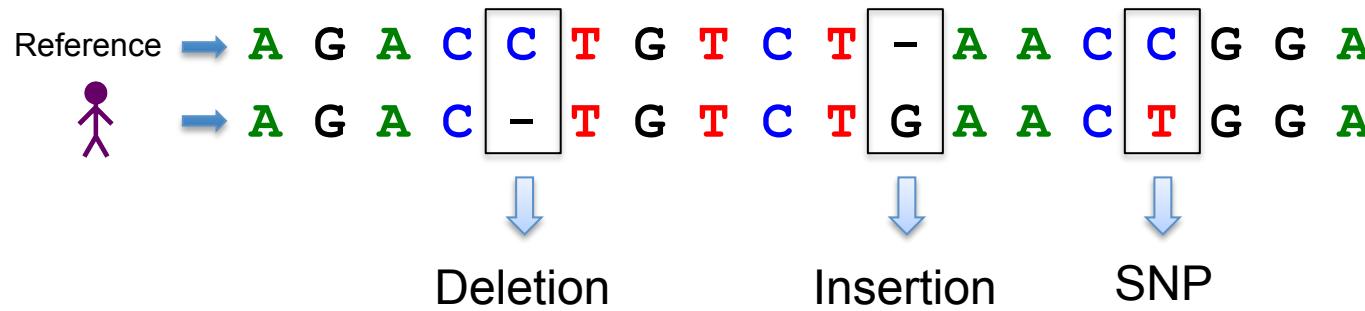
Two human genomes differ by one base in every 1000 on average

- We study genomic variation to understand:
 - Why some individuals are more susceptible to disease
 - Why some individual have an adverse response to drugs or environmental factors
 - Population history
- There are over **150 million** short human variants in dbSNP (v144)
- There are over **3 million** longer (structural) variants in the Database of Genomic Variants archive (DGVa)

Variation types - Short variants

Short scale: one or few nucleotides

- Single nucleotide polymorphism (SNP)
- Small insertions and/or deletions (DIPs or indels)



Variation types - Long variants

Large scale: in chromosomal structure (structural variation)

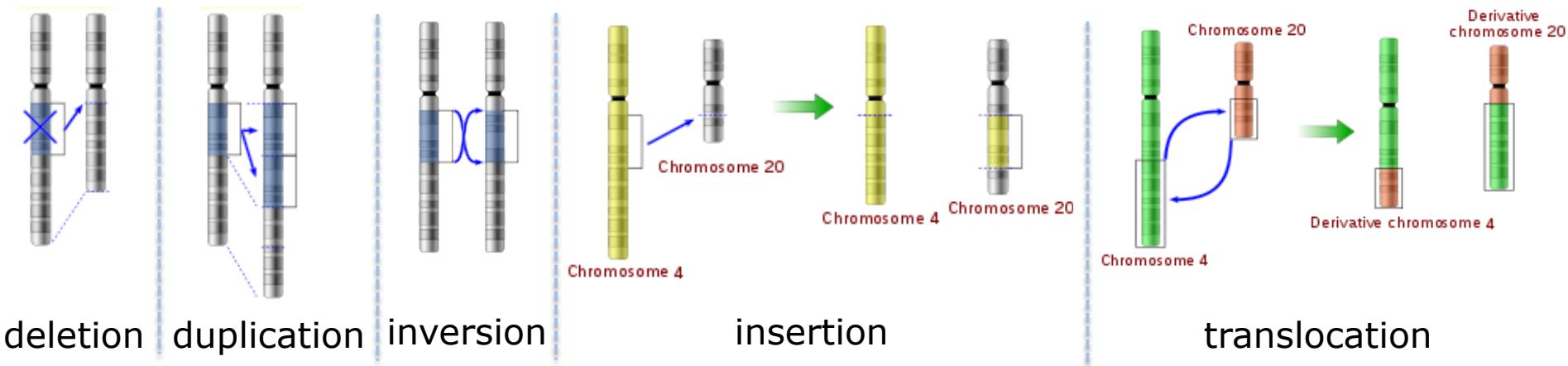
- Copy number variations (CNVs): sequence repeated 'n' times in an individual

Reference (4): 

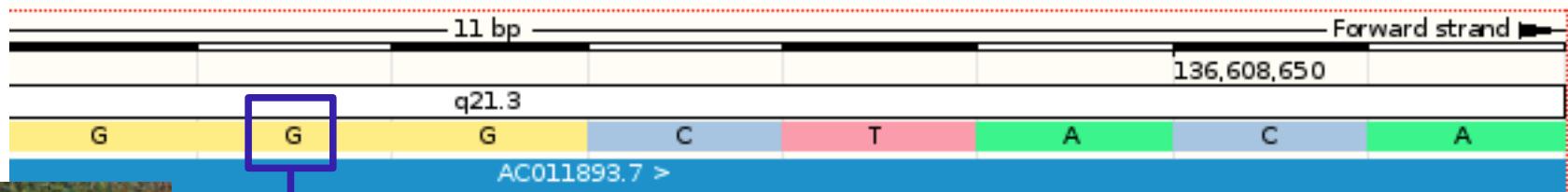
Copy gain (6): 

Copy loss (3): 

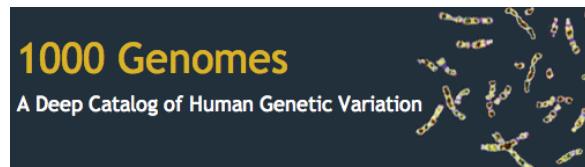
- Large structural variants:



Variation



GAGCTACA
GAGCTACA
GGGCTACA
GGGCTACA



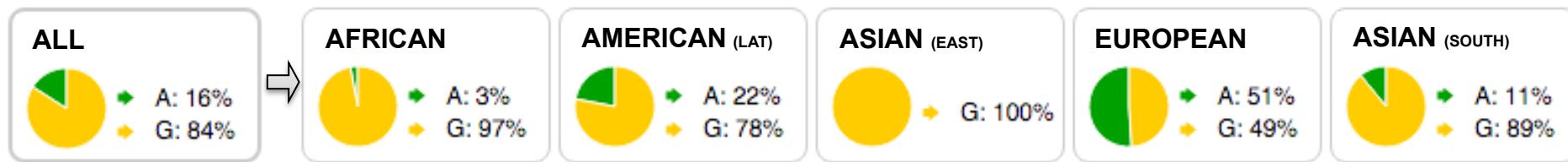
Variation

Lactose Intolerance

rs4988235: a SNP near the LCT gene controls whether lactase enzyme is turned on or off as a person grows older.

Alleles: G/A

1000 Genomes Project Phase 3 allele frequencies



Genotype	What it means
GG	Likely to be lactose intolerant.
GA	Likely to be tolerant due to lactase persistence.
AA	



Human (GRCh38.p5)

Location: 9:22,125,004-22,126,004

Variant: rs1333049

Search all species...

Variant displays

- Explore this variant
- Genomic context
 - Genes and regulation
 - Flanking sequence
- Population genetics
- Sample genotypes
- Linkage disequilibrium
- Phenotype Data
- Phylogenetic Context
- Citations
- External Data
 - LOVD

Configure this page

Add your data

Export data

Share this page

Bookmark this page

rs1333049 SNP

Original source

Variants (including SNPs and indels) imported from dbSNP (release 144) | [View in dbSNP](#)

Alleles

G/C | Ancestral: C | Ambiguity code: S | MAF: 0.42 (C)

Location

Chromosome 9:22125504 (forward strand) | [View in location tab](#)

Most severe consequence

Downstream gene variant | [See all predicted consequences](#) [Genes and regulation]

Evidence status

[Archive dbSNP rs59077428](#), [rs58844516](#), [rs17761501](#)

g:g.22125504G>C

HGVS name

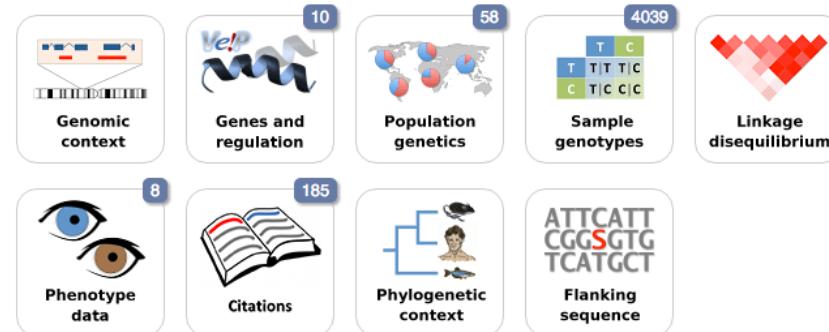
This variant has assays on 7 chips - [Show](#)

Genotyping chips

This variant overlaps 10 transcripts, has 4039 sample genotypes, is associated with 8 phenotypes and is mentioned in 185 citations.

About this variant

Explore this variant



Using the website

- Video: [Browsing SNPs and CNVs in Ensembl](#)
- Video: [Clip: Genome Variation](#)
- Video: [BioMart: Variation IDs to HGNC Symbols](#)
- Exercise: [Genomes and SNPs in Malaria](#)

Analysing your data



Test your own variants with the Variant Effect Predictor

Programmatic access

- Tutorial: [Accessing variation data with the Variation API](#)

Reference materials

- [Ensembl variation documentation portal](#)
- [Ensembl variation data description](#)
- [Variation Quick Reference card](#)

http://www.ensembl.org/Homo_sapiens/Variation/Explore?v=rs1333049

Ensembl Variation - Roles

Data:

- Build the variation databases: import variations from [different sources](#)
- Quality control ([quality filtering](#), [evidence summary](#))
- Data annotation: e.g. phenotype data, ancestral alleles, ...
- Calculate [consequences](#) of the variants on the transcripts
- [SIFT](#) and [PolyPhen](#) scores

API:

- Build and maintain the Variation API

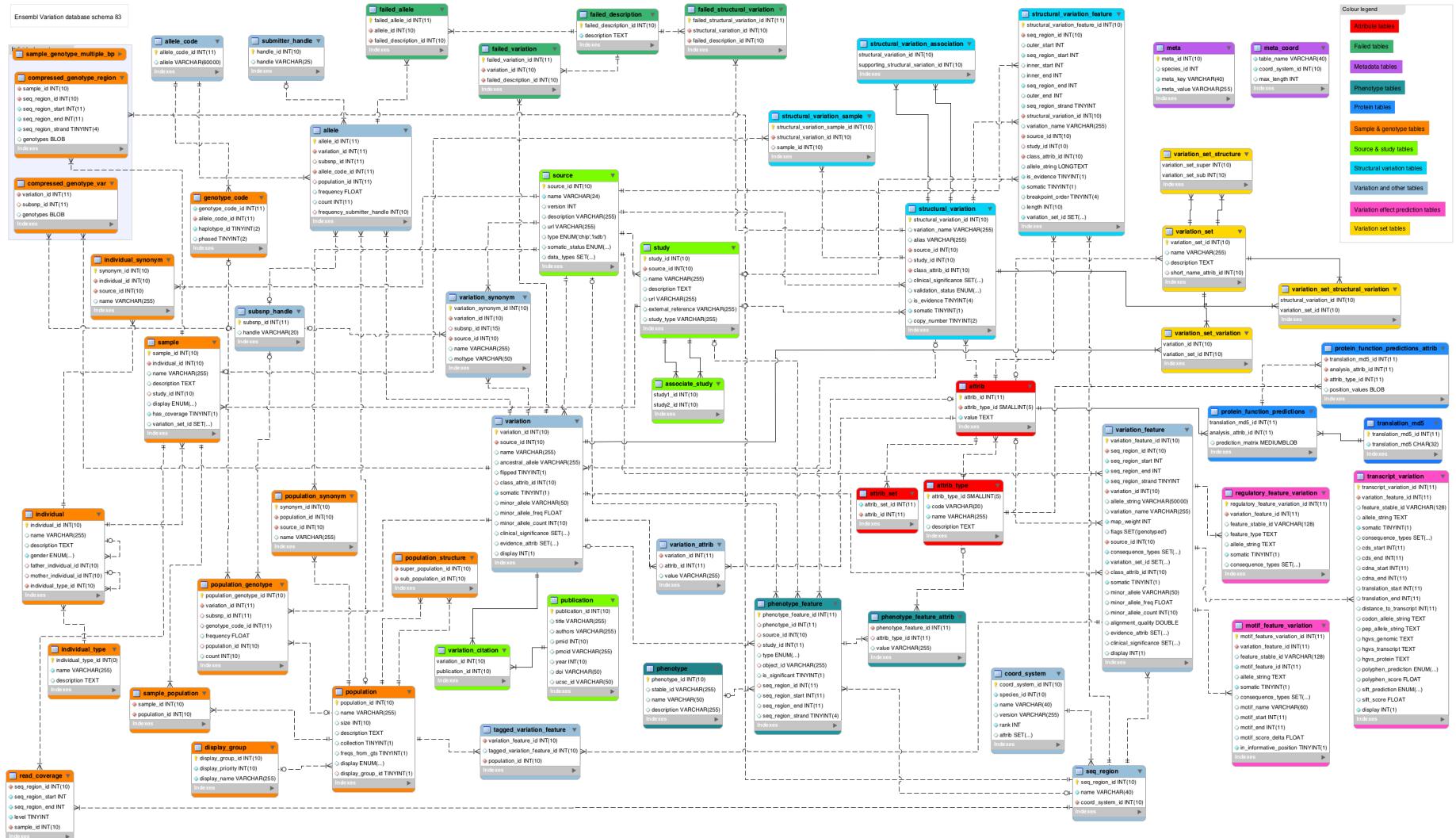
Web:

- Build and maintain part of the Variation Web API



22

Variation database schema (MySQL)



→ http://www.ensembl.org/info/docs/variation/variation_schema.html



eEnsembl BLAST/BLAST | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Human (GRCh37) | Location: 2:136,608,146-136,609,146 | Variation: rs4988235

rs4988235 SNP

Original source: Variants (including SNPs and indels) imported from dbSNP (release 137) | [View in dbSNP](#)

Reference/Alternative: G/A | Ancestral: G | Ambiguity code: R1 MAF: 0.23 (A)

Chromosome 2:136,608,146 (forward strand) | [View in location tab](#)

Co-located with HGMD-PUBLIC: C802420

1000 Genomes, HepMap, Cited, Frequency, Multiple observations

Synonyms: None currently in the database

This variation has 4 HGVS names - click the plus to show:

2:g.136608646G>A
ENST00000483902.1:n.1917<326C>T
ENST00000483902.1:n.343<326C>T
ENST00000492091.1:n.343<326C>T

This variation has assays on 4 chips - click the plus to show:

Explore this variation

- Genomic context
- Genes and regulation
- Population genetics
- Individual genotypes
- Linkage disequilibrium
- Phenotype data
- Phylogenetic context
- Flanking sequence

Configure this page

[Add your data](#)

[Export data](#)

[Bookmark this page](#)

[Share this page](#)

```

use Bio::EnsEMBL::Registry;
use Bio::EnsEMBL::MappedSliceContainer;
use Bio::EnsEMBL::DBSQL::StrainSliceAdaptor;
use Bio::EnsEMBL::DBSQL::AssemblySliceAdaptor;

# get registry
my $reg = 'Bio::EnsEMBL::Registry';

my $registry_file = '.ensembl.registry';
$reg->load_all($registry_file);
#$reg->load_registry_from_db(-host => 'ensembldb.ensembl.org', -user => 'anonymous');

my $sa = $reg->get_adaptor("human", "core", "slice");

my $slice;

if(scalar @ARGV) {
    $slice = $sa->fetch_by_region('chromosome', $ARGV[0], $ARGV[1], $ARGV[2]); # simon's long
}

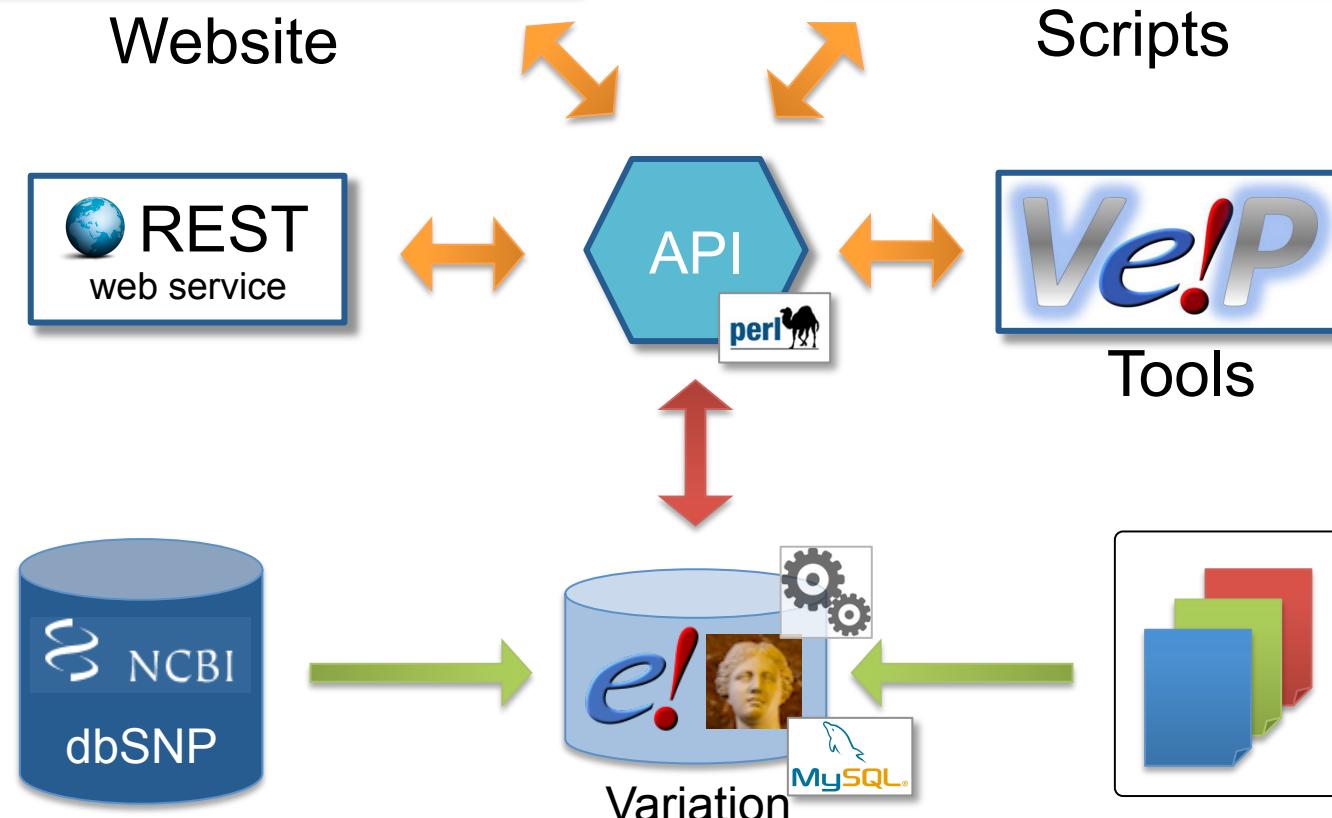
else {
    $slice = $sa->fetch_by_region('chromosome', 21, 34698530, 34698570);
    #$slice = $sa->fetch_by_region('chromosome', 13, 35016110, 35016140);
}

# create a new mapped slice container
my $msc = Bio::EnsEMBL::MappedSliceContainer->new(-SLICE => $slice, -EXPANDED => 1);

# create a new strain slice adaptor and attach it to the mapped slice container
my $ssa = Bio::EnsEMBL::DBSQL::StrainSliceAdaptor->new($sa->db);
$msc->set_StrainSliceAdaptor($ssa);

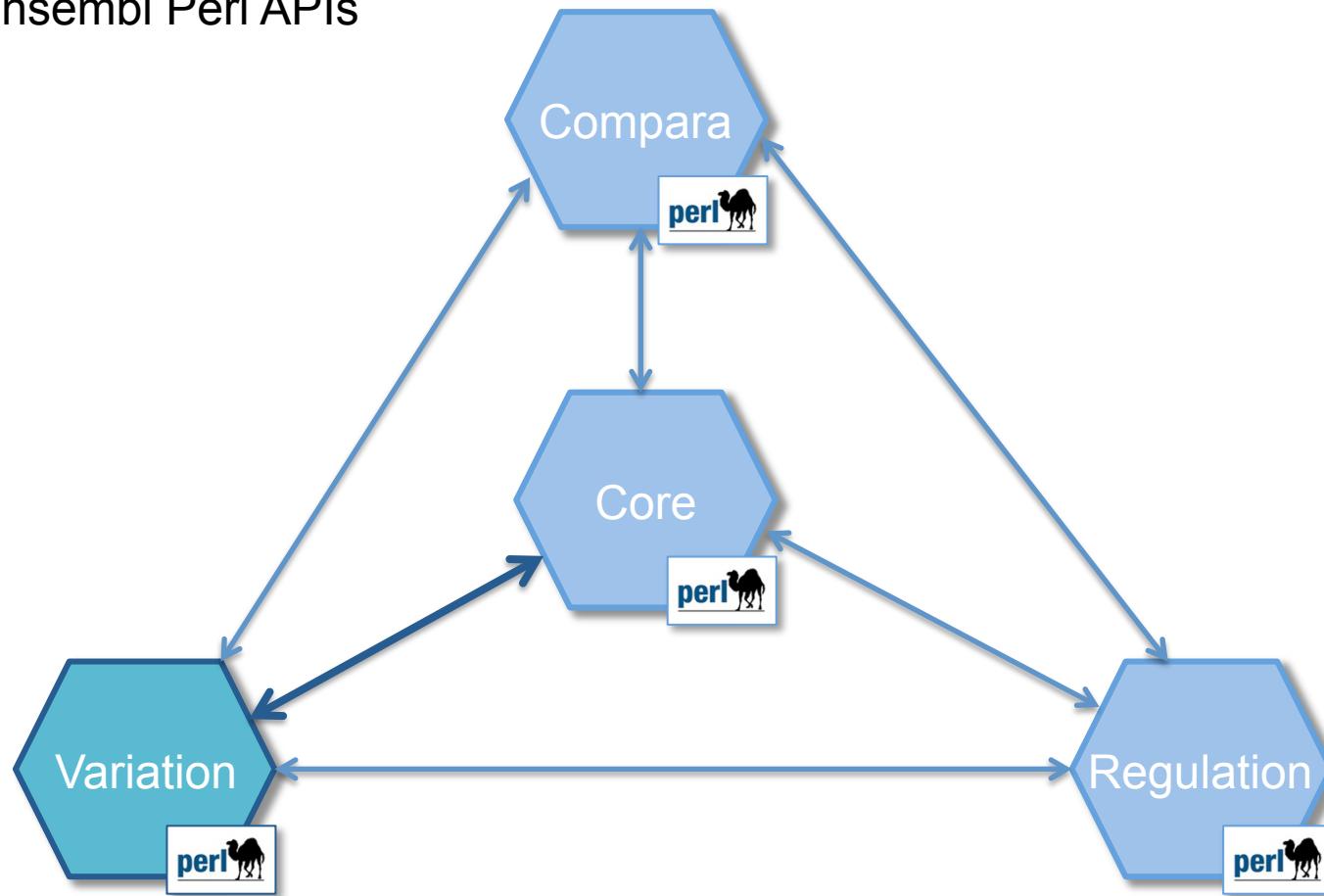
# now attach strains
$msc->attach_StrainSlice('Watson');
$msc->attach_StrainSlice('Venter');

```



About Ensembl API

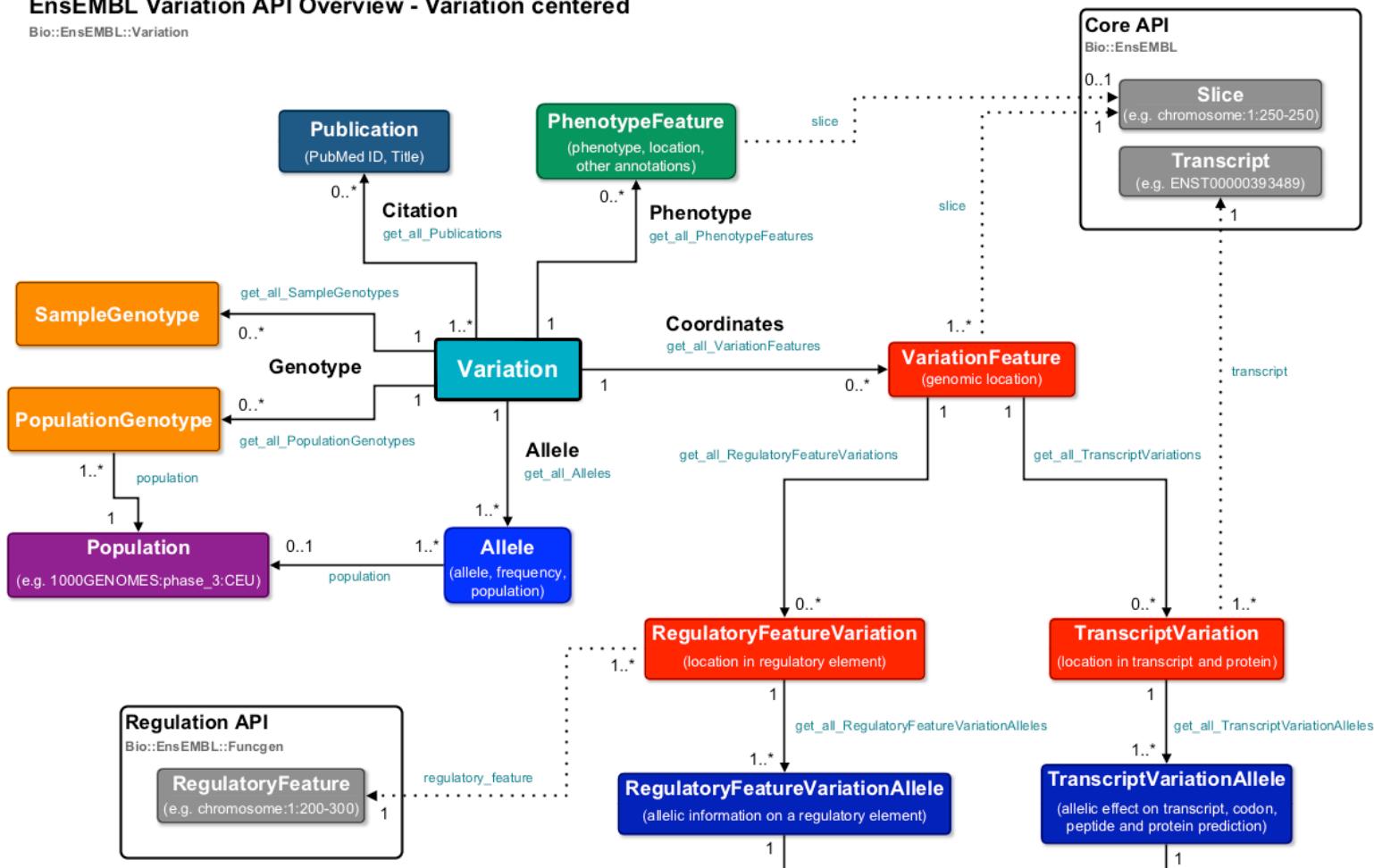
4 main Ensembl Perl APIs



Simple objects diagram

EnsEMBL Variation API Overview - Variation centered

Bio::EnsEMBL::Variation



Legend

- Central class
- Feature class
- Class from a different Ensembl API

text Method to get the object(s)
 → Get object(s) from the same API
 → Get object(s) from a different API

http://www.ensembl.org/info/docs/api/variation/variation_API_diagram.html

Objects and adaptors

- The API deals with objects representing database entities
- Adaptors are “**factories**” for generating objects
 - Adaptors are retrieved from the Registry

```
use Bio::EnsEMBL::Registry;

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

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

my $va = $reg->get_adaptor('human', 'variation', 'variation');

my $variation = $va->fetch_by_name('rs334');
```

The diagram consists of three colored arrows pointing upwards from labels to their corresponding arguments in the code:

- A red arrow points from the label "species" to the argument "human".
- A blue arrow points from the label "group" to the argument "variation".
- A green arrow points from the label "object name" to the argument "rs334".

Object creation

- Using adaptors
 - **Fetch** object(s) according to some property e.g. name, location
 - “fetch_all_...” -> returns a list reference of items
 - “fetch_by_...” -> usually returns only 1 item
 - Check documentation which methods the adaptor provides
- Using API objects: e.g. **Variation**
 - **Get** other object(s) from an API object
 - e.g. \$variation->**get_all_Alleles()**
returns a list reference of Ensembl **Allele** objects
 - Usually the object is written with a upper case in the method
 - e.g. **get_all_Alleles()**

Key data/objects

Variation object

Represents a short difference between at least two genomic sequences

- Basic unit in Ensembl Variation database
- Retrieve using variation adaptor
- May have an evidence class, ancestral allele, clinical significance
- Key attributes:

Attribute	Example value(s)	Method(s)
Variant name	rs1333049, COSM679126	<code>\$v->name()</code>
Source	dbSNP, COSMIC	<code>\$v->source_name()</code>
Class	SNP, insertion, deletion	<code>\$v->var_class()</code>
Minor allele	Allele: C Frequency: 0.43	<code>\$v->minor_allele()</code> <code>\$v->minor_allele_frequency()</code>
Supporting evidence	Array: [Frequency, Multiple_observations,etc...]	<code>\$v->get_all_evidence_values()</code>

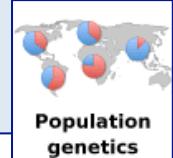
Exercise 1

- Retrieve the following information:
 - Variation class
 - Source
- for the following variations in human:
 - rs1333049
 - rs56385407
 - COSM998
 - CI003207

Hint: You will need to use the **VariationAdaptor** and the method “**fetch_by_name**” to retrieve the Variation objects

- Variation API documentation:
<http://www.ensembl.org/info/docs/Doxygen/variation-api/index.html>
- Variation tutorial:
http://www.ensembl.org/info/docs/api/variation/variation_tutorial.html

Allele information

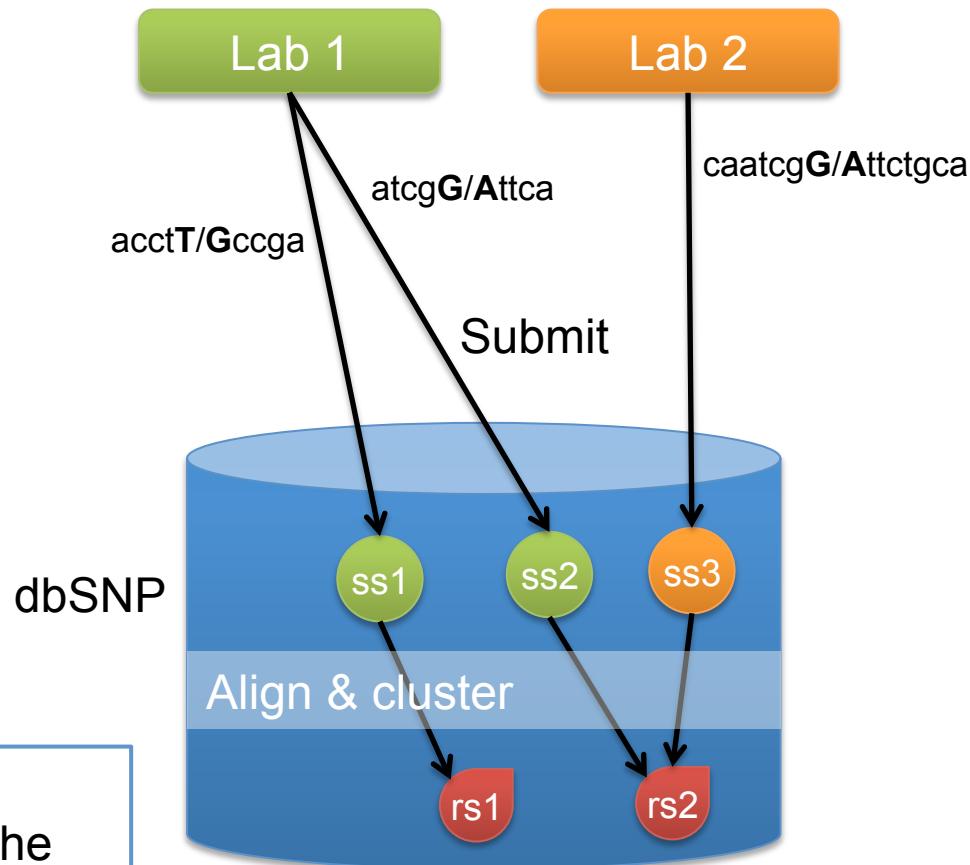


A variant at a specific location can have two or more alleles.
Usually a reference allele and at least 1 alternative allele.

Information held against an allele includes:

- Base(s) observed
- Source population
- Submitter information (dbSNP ‘handle’)
- Frequency in an assayed population

Frequencies for large studies, such as the 1000 Genomes project are calculated on the fly from VCF files



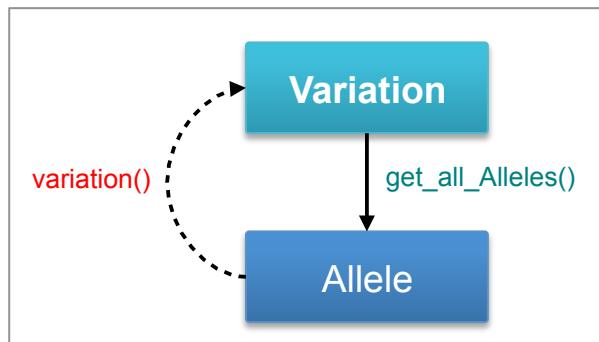
Allele objects

- Represents an allele observed at a variation and usually in a population
- Can be retrieved from **Variation** objects, as well as the **Allele adaptor**

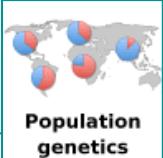
Attribute	Example value(s)	Method(s)	Comment(s)
Allele	G, AC	<code>\$a->allele()</code>	
Frequency	0.15, 1	<code>\$a->frequency()</code>	Not always defined
Population	Bio::EnsEMBL::Variation::Population	<code>\$a->population()</code>	Returns an object if defined
Variation	Bio::EnsEMBL::Variation::Variation	<code>\$a->variation()</code>	Returns an object
Submitter id	TSC_CSHL	<code>\$a->subsnp_handle()</code>	Not always defined

Note: most API objects have a method to retrieve reference to “parent” object
 e.g. to retrieve variation object from allele object:

```
$allele->variation()
```



Exercise 2



- For SNP rs1333049 in human, retrieve the following for each of its alleles:
 - Allele
 - Frequency *
 - Population name *
 - Submitter name in dbSNP (“handle”) *

* if exist

Note:

Phase 3 data from the 1000 Genomes Project is accessed directly from VCF files. It is not default API behaviour to read VCF's, so to see frequencies for this project you will need to set the 'use_vcf' option on the database adaptor to 1.

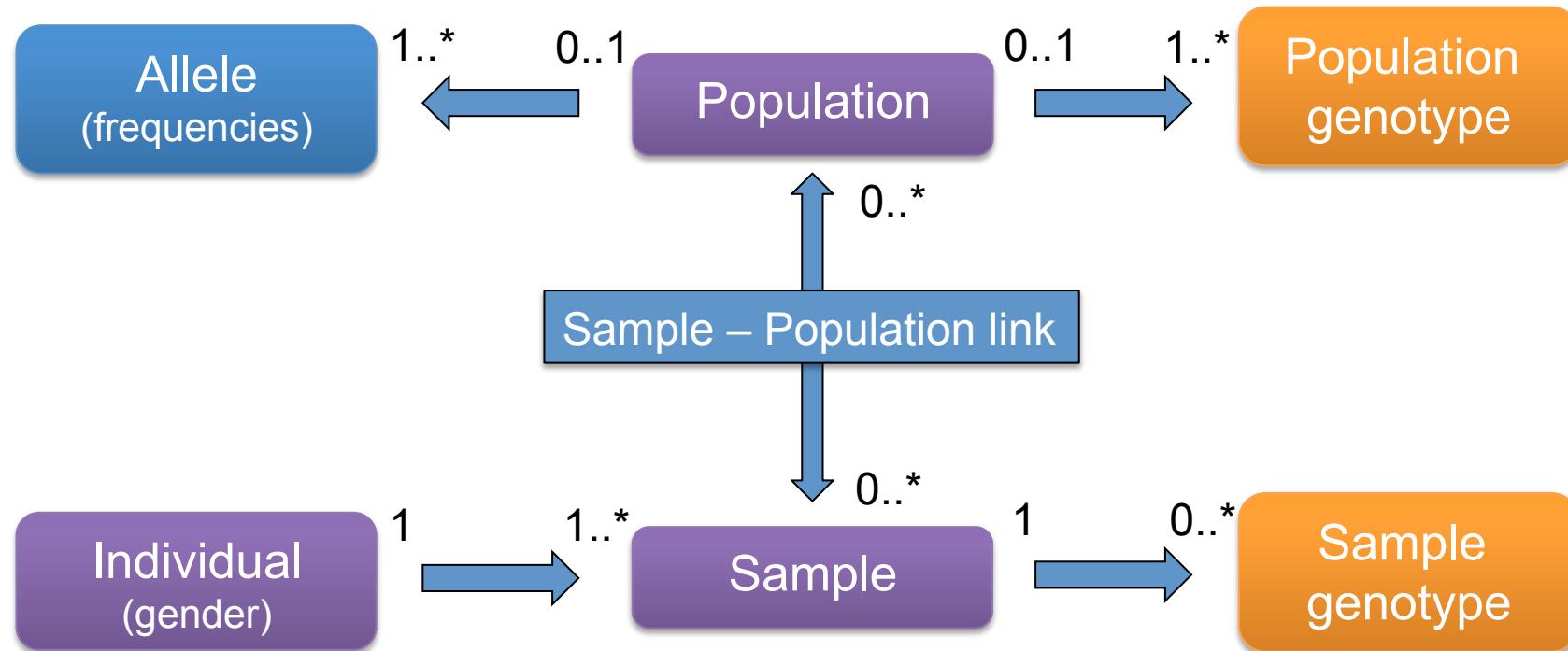
```
$va->db->use_vcf(1);
```

Hint:

You may need to test whether the Population object returned is empty to avoid a script error.

Sample, Individuals & Populations

A population is a group of samples used in a study, rather than a geographical grouping



Experimental results are attached to samples, which are taken from individuals; this supports results from multiple samples from the same individual

Sample Genotypes

Sample genotypes show which combination of the possible alleles a specific sample from a specific individual has at the site of a specific variant.

We import genotype data from dbSNP and large scale sequencing project such as the 1000 Genomes project, the WTSI Mouse Genomes Project and the NextGen Project.

Genotype data is large, so it is stored either in VCF files or in a compressed format in the Ensembl databases. It cannot be accessed using direct SQL



Sample Genotype object

Represents an instance of the genotype of a variant in a specific sample from an individual

Can be retrieved from **Variation** object, as well as the
SampleGenotype adaptor (Bio::EnsEMBL::Variation::DBSQL::SampleGenotypeAdaptor)

Attribute	Example value(s)	Method(s)
Genotype	List reference [A,C]	\$sg->genotype()
Genotype string	A C	\$sg->genotype_string()
Variation object	Bio::EnsEMBL::Variation::Variation	\$sg->variation()
Sample object	Bio::EnsEMBL::Variation::Sample	\$sg->sample()
Individual object	Bio::EnsEMBL::Variation::Individual	\$sg->sample->individual()

Exercise 3

- Fetch all available individual genotypes for the human variation rs1333049.
- Extract:
 - The genotype string
 - The name of the sample
 - The gender of the individual

Note:

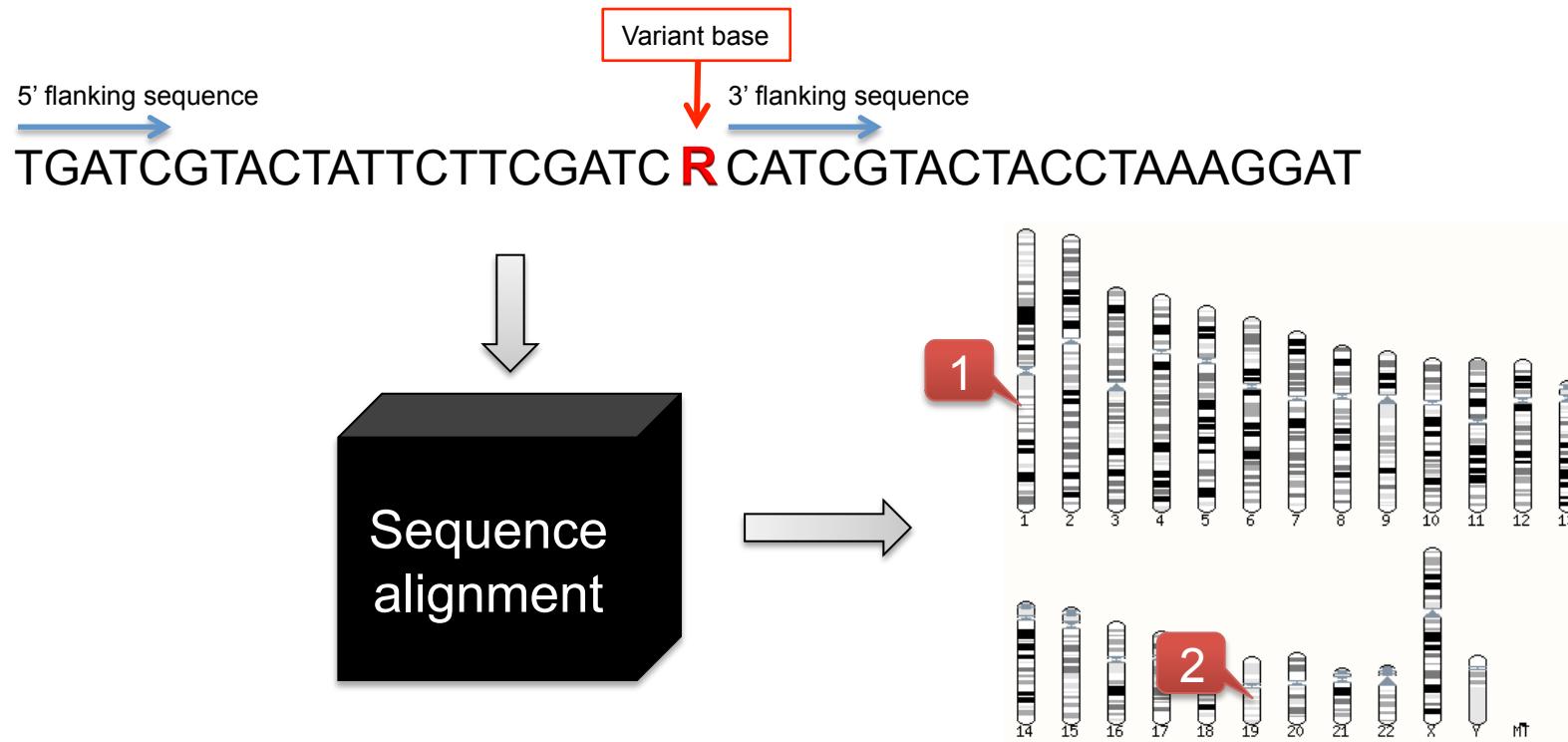
Phase 3 Data from the 1000 Genomes Project is accessed directly from VCF files. It is not default API behaviour to read VCF's, so to extract genotypes from this project you will need to set the '**use_vcf**' option on the database adaptor to 1:

```
$va->db->use_vcf(1);
```

Variation mapping

Variations are mapped to the genomic reference sequence using their flanking sequences

- A variation may have multiple mappings on the current reference
- A variation may not map reliably to the current reference



Variation feature object

- An instance of a variation mapping to the genome
- Can be retrieved from **Slice** (Core) and **Variation** objects, as well as **VariationFeature adaptor**

Attribute	Example value(s)	Method(s)	Comment(s)
Allele string	A/G, -/C	<code>\$vf->allele_string()</code>	
Chromosome	15, X	<code>\$vf->seq_region_name()</code>	
Coordinates	103019234	<code>\$vf->start()</code> <code>\$vf->end()</code> <code>\$vf->seq_region_start()</code> <code>\$vf->seq_region_end()</code>	<p>} slice relative</p> <p>} chromosome relative</p>
Slice object	Bio::EnsEMBL::Slice	<code>\$vf->slice()</code>	Returns a Core API object

Note: Allele strings are reported as Reference/Alternative where possible:

A/G → “reference/alternative”

Exercise 4

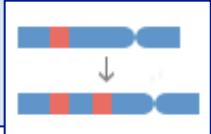
- Retrieve all variations in human located on chromosome 13 from 48987260 to 48990000 and get:
 - Variation name
 - Alleles (“allele_string”)
 - Location (e.g. Chromosome:Start-End)
- Extra:

Find the the genomic location(s) of the following human variants:

 - rs7107418
 - rs671
 - rs17646946
 - rs4988235

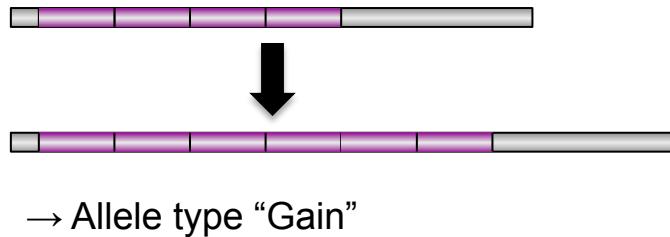
Hint: usually when you have to search something by location, you need to use the **Slice** object ([Core API](#)) first

Structural variation

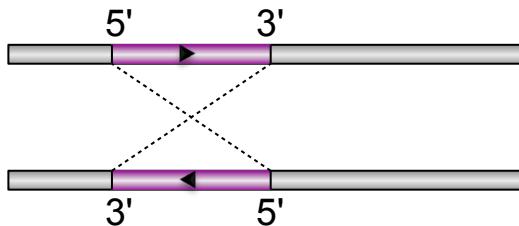


- Usually represents a large segment of variable DNA sequence (> 1kb)
- Different processes give rise to different types of structural variations:

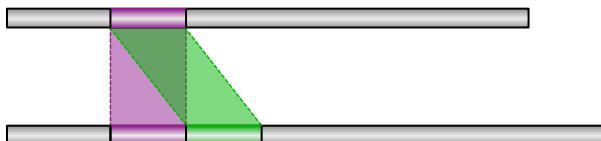
Copy number variation (CNV)



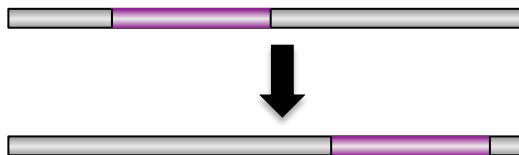
Inversion



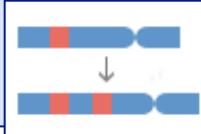
Duplication



Translocation



Example of structural variation



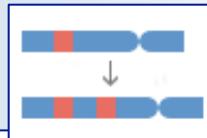
Structural variant: [nsv917561](#)

Variation class	CNV (SO:0001019)
Source	DGVa – Database of Genomic Variants Archive
Study	nstd75 – “International Standards for Cytogenomic Arrays Consortium (prenatal dataset).” PMID:21844811, PMID:20466091
Location	Chromosome 1:146987841-148359881 (forward strand)

Supporting evidence

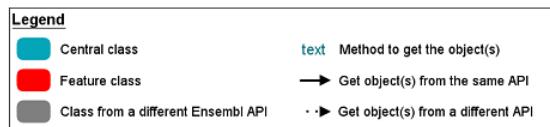
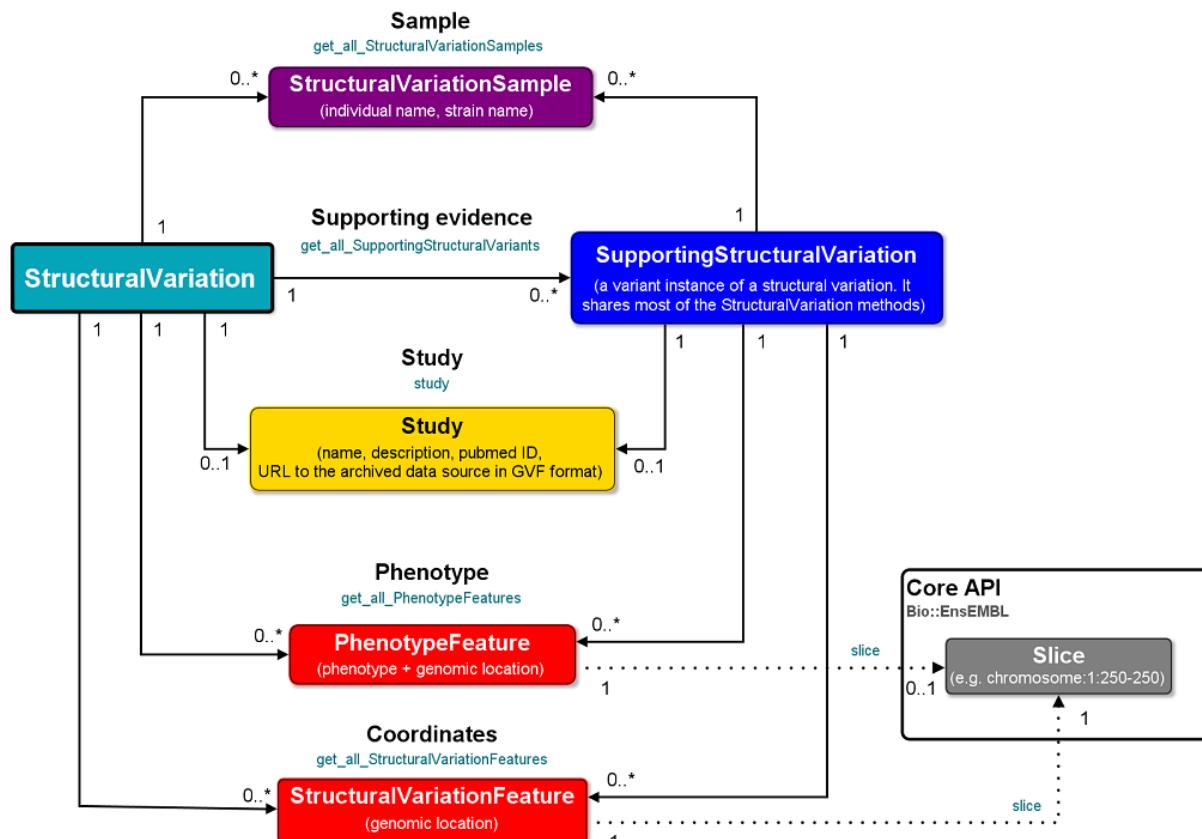
Supporting evidence	Chr:bp (strand)	Allele type	Clinical significance	Individual name
nssv1605931	1:146987841-148359881 (+)	Gain	pathogenic	ISCA_ID_pn_1622
nssv1606072	1:146987841-148359881 (+)	Gain	-	ISCA_ID_pn_1746

Simple objects diagram



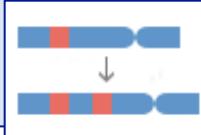
EnsEMBL Variation API Overview - StructuralVariation centered

Bio-EvoEMBI · Variation



 http://www.ensembl.org/info/docs/api/variation_structuralvariation_API_diagram.html

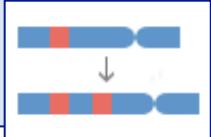
Structural variation object



- Similar to the **Variation** object, but represents longer variants
- Data from:
 - DGVa (with study and supporting evidences)
 - CNV probes from Illumina and Affymetrix

Attribute	Example value(s)	Method(s)	Comment(s)
Structural variation name	esv214236	<code>\$sv->variation_name()</code>	
Class	CNV	<code>\$sv->var_class()</code>	
Study	Bio::EnsEMBL::Variation::Study	<code>\$sv->study()</code>	Returns an object if defined
Supporting evidence(s)	Reference list of objects Bio::EnsEMBL::Variation::SupportingStructuralVariation	<code>\$sv->get_all_SupportingStructuralVariants()</code>	

Supporting structural variation



Supporting evidence <=> Allele for a Variation (i.e. experimental data)

StructuralVariation

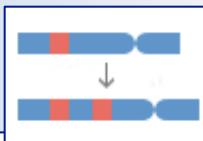
Region of the genome that a submitter has defined as containing structural variations.

get_all_SupportingStructuralVariants

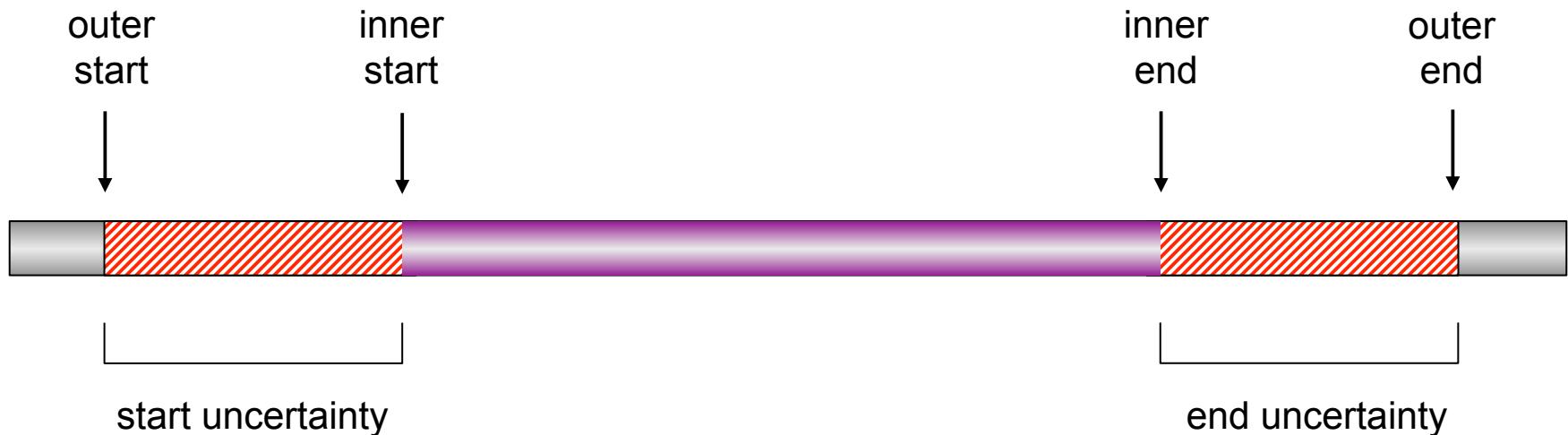
SupportingStructuralVariation

Actual **variant calls** that were made within a study.

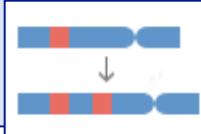
Structural variation coordinates



- Usually a structural variation has only a start and an end
- But sometimes, the breakpoint locations cannot be determined precisely



Structural variation feature object



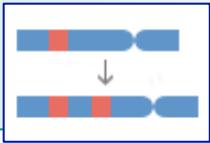
- Similar to **VariationFeature** but represent structural variations
- Can be retrieved from a **Slice** (Core) as well as the **StructuralVariationFeature adaptor** and the **StructuralVariation** object
- No allele string, alleles or associated variation
- Only coordinates, class and the list of the supporting evidences

Attribute	Example value(s)	Method(s)
Structural variation name	esv214236	<code>\$svf->variation_name()</code>
Chromosome	15, X	<code>\$svf->seq_region_name()</code>
Coordinates		<code>\$svf->start()</code> <code>\$svf->end()</code> <code>\$svf->seq_region_start()</code> <code>\$svf->seq_region_end()</code>
		<code>\$svf->outer_start()</code> <code>\$svf->inner_start()</code> <code>\$svf->inner_end()</code> <code>\$svf->outer_end()</code>
Class	CNV	<code>\$svf->var_class()</code>

“Classic”
coordinates

SV specific
coordinates

Exercise 5



- Fetch the following information for the structural variation [esv275264](#) in Human :
 - Structural variation class
 - Study name and description
 - Coordinates
- Fetch the names and the classes (sequence ontology term) of its supporting structural variations.

Additional information/annotation

Variation Sets

Variation sets are arbitrary grouping of variations

- Useful for limiting scripts to specific subsets of data
- May group data from multiple sources
- May be arranged in a parent set/subset formation

Parent set	Set
All phenotype-associated variants	HGMD-PUBLIC variants, ClinVar , OMIM, NHGRI-EBI GWAS catalog...
1000 Genomes 3 - All	1000 Genomes 3 - AFR – common, 1000 Genomes 3 – EUR...
Genotyping chip variants	Affy GeneChip 500K, Illumina_HumanOmni2.5 Illumina_Cardio-Metabo_Chip
	ESP_6500

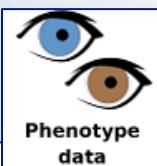
http://www.ensembl.org/info/genome/variation/data_description.html#variation_sets

Variation Sets

- Arbitrary collections of **Variations** or **Structural variations**
- Useful to limit your scripts to important subsets
- Sometimes contain millions of **Variations**, in which case you can fetch an **Iterator** instead of a list

Attribute	Example value(s)	Method(s)
Description	Variations called by the 1000 Genomes Project...	<code>\$vs->description()</code>
Members	List of Variations Iterator over Variations	<code>\$vs->get_all_Variations()</code> <code>\$vs->get_Variation_Iterator()</code>
Sub/Super sets	Bio::EnsEMBL::Variation::VariationSet	<code>\$vs->get_all_sub_VariationSets()</code> <code>\$vs->get_all_super_VariationSets()</code>

Phenotype data



ClinVar

european
genome-phenome
archive



IMPC
INTERNATIONAL MOUSE
PHENOTYPING CONSORTIUM

GIANT
CONSORTIUM



OMIM®

DGVA^{rchive}

COSMIC
Catalogue of somatic mutations in cancer

DECIPHER
GRCh37

orphanet

LOVD
Leiden Open Variation Database

dbGaP
GENOTYPES and PHENOTYPES

Animal QTLdb

ZFIN

MAGIC
Metabolism of Glucose and Insulin-related Traits Consortium

UniProt



OMIA - ONLINE MENDELIAN INHERITANCE IN ANIMALS



Variant

Gene

Structural Variant

QTL

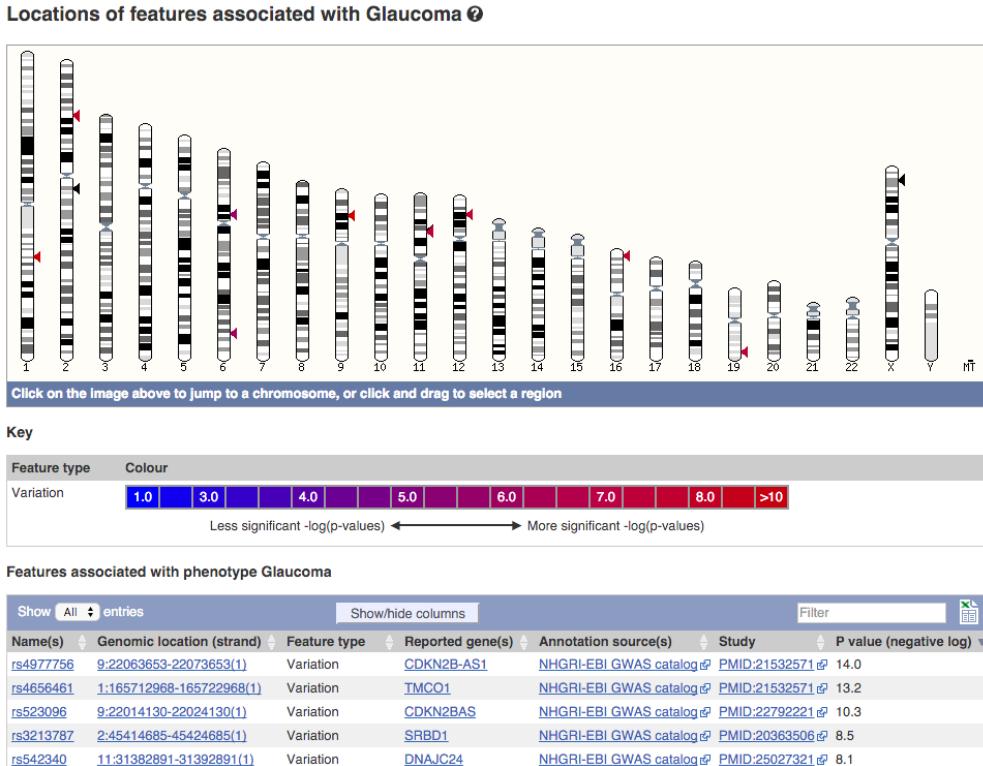
Phenotype data

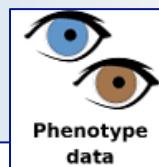
Phenotypes are mapped to genomic locations via the variations or genes they have reported associations with to create a **PhenotypeFeature**

We are collaborating with groups working on phenotype ontologies, but the data current held is as it appears in the source database.

Attributes types held for phenotype features include:

- Clinical significance from ClinVar
- Reported genes
- Most associated risk allele
- P-value
- Inheritance type
- Odds ratio
- Beta coefficient





Phenotype Feature object

Represents an association between a phenotype and a genomic feature

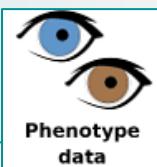
- Can be retrieved from a **Slice** (Core) or a **Variation** object as well as its **PhenotypeFeature adaptor**
- Not all the data has the same attributes populated
- Phenotype + genomic location

Attribute	Example value(s)	Method(s)
Phenotype object	Bio::EnsEMBL::Variation::Phenotype	\$pf-> phenotype () *returns object
Reported gene(s)	CDKN2BAS, NOTCH2	\$pf-> associated_gene ()
p-value	6e-07	\$pf-> p_value ()
Risk allele	C	\$pf-> risk_allele ()
Chromosome	15, X	\$vf-> seq_region_name ()
Coordinates	103019234	\$vf-> start () \$vf-> end () \$vf-> seq_region_start () \$vf-> seq_region_end ()

To get the phenotype name:

\$pf->**phenotype->description** ()

Exercise 6



Find all phenotypes associated with the human variants **rs6495122** and **rs2470893** and extract:

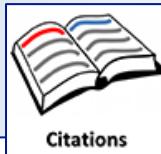
- The phenotype description
- The source name
- The p-value for the association
- The reported risk allele
- Any reported gene(s)

Extra

For each variant find:

- The minor allele and its frequency (this is calculated from the full set of 1000 Genomes samples)
- The ancestral allele (this is calculated by aligning 8 primate assemblies)

Variation citation data



We import variation citation data from:

- **dbSNP**
 - Publication information is submitted with the variant submission.
- **EuropePMC**
 - Publications for which the full text is freely available in PubMed Central are mined for refSNP identifiers.
- **UCSC Genocoding Project**
 - Publications for which the text is freely available of those from publishers who have agreed to data mining are mined for refSNP identifiers.

In Ensembl 83, **9 species** have [citation data](#).

For human we have:

- 156,000 distinct variants associated with at least 1 publication
- 53,000 distinct publications

Publication object



- Represents an article in a journal
- Can be retrieved from a **Variation** object as well as the **Publication adaptor**
- The type of information available will depend on publisher and data source

Attribute	Example value(s)	Method(s)
Cited variants	List reference of Bio::EnsEMBL::Variation::Variation objects	\$p-> variations() <i>*returns listref of objects</i>
Title	Mitochondrial acetylation and diseases of aging.	\$p-> title()
Authors	Wagner GR, Payne RM.	\$p-> authors()
Year	2011	\$p-> year()
PubMed Identifier	21437190	\$p-> pmid()

Exercise 7



Obtain Variation objects for the human variations:

- rs3730070
- rs11765954
- rs671

Find all articles we hold which cite these variants and report :

- The Pubmed ID
- The year of publication
- The title

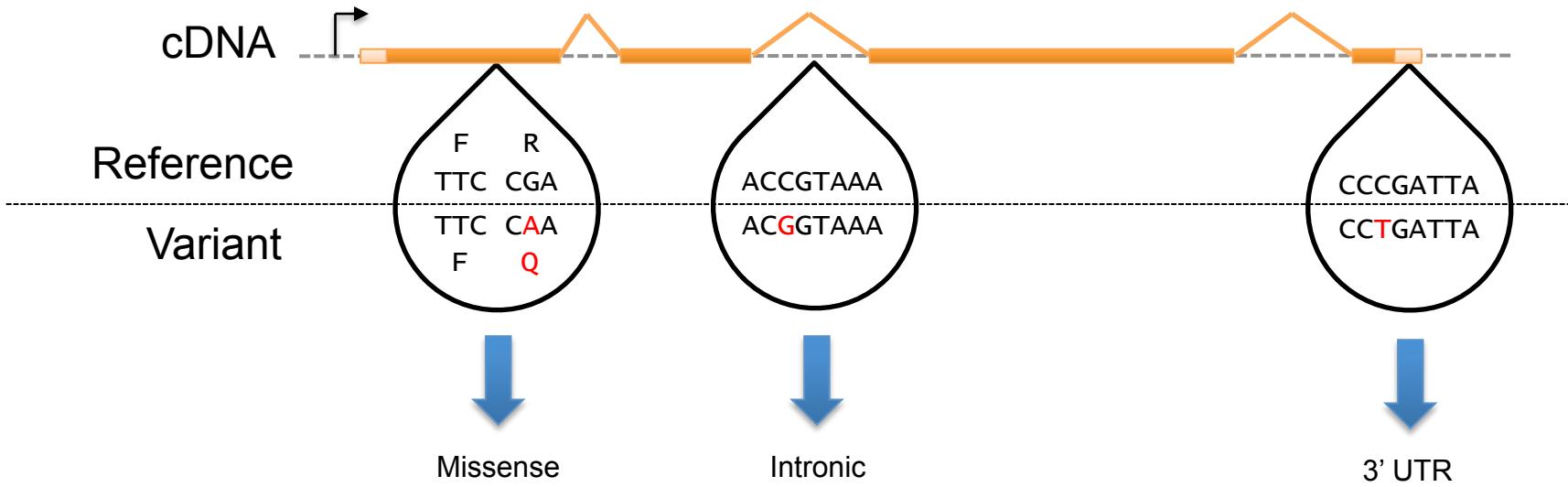
Extra:

Find all variants cited in the publication with the PMID:18187665

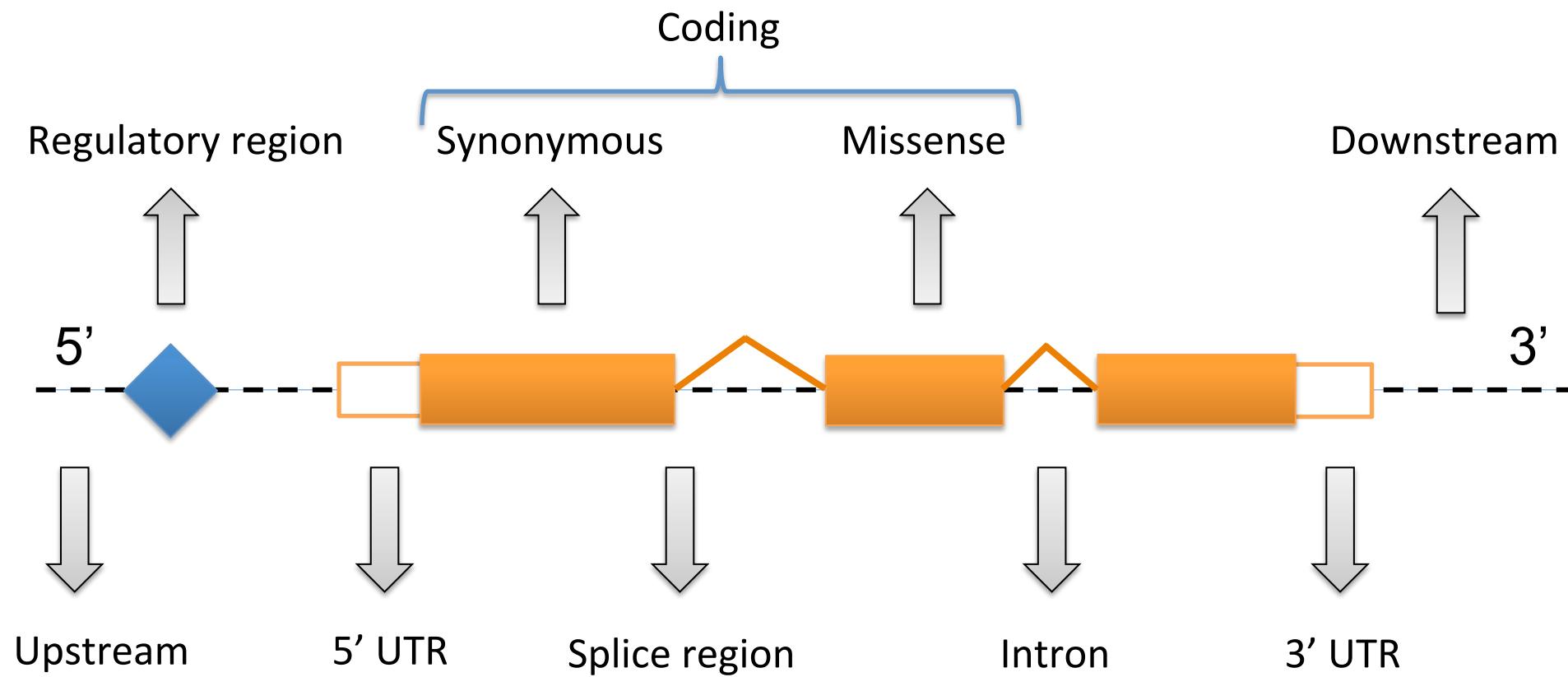
Analysis

Variation consequences

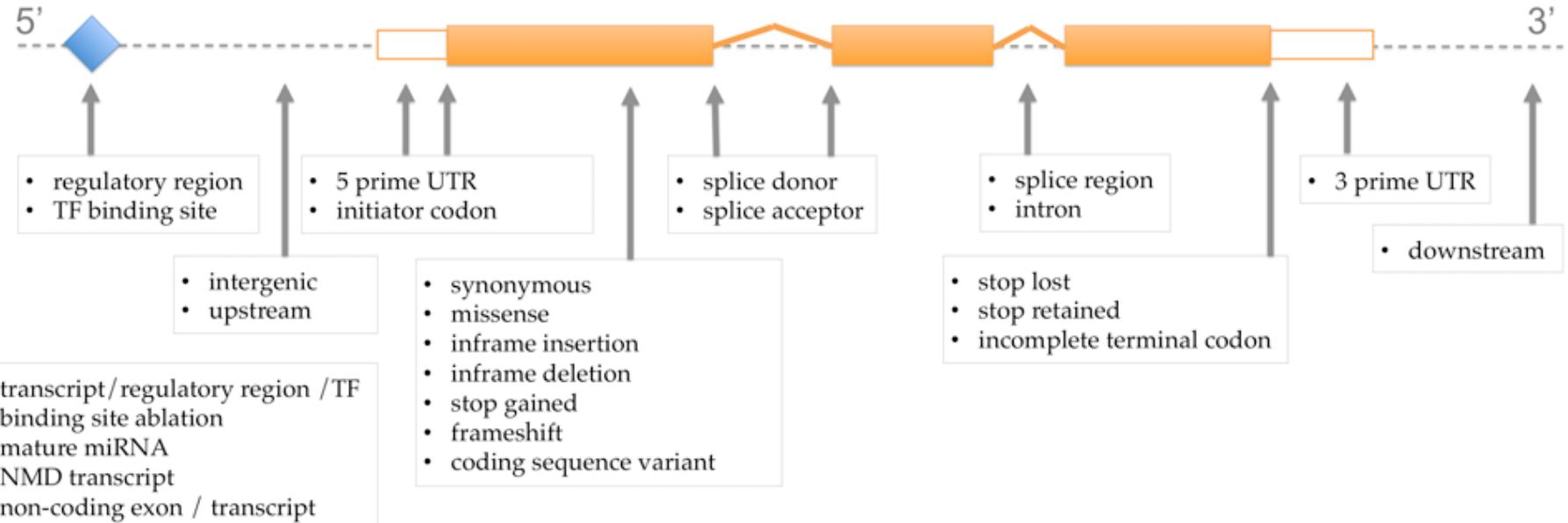
Defines consequence of a variation in relation to the transcript structure it overlaps



Variation consequences



Variation consequences



Sequence ontology provides controlled vocabulary to describe consequences

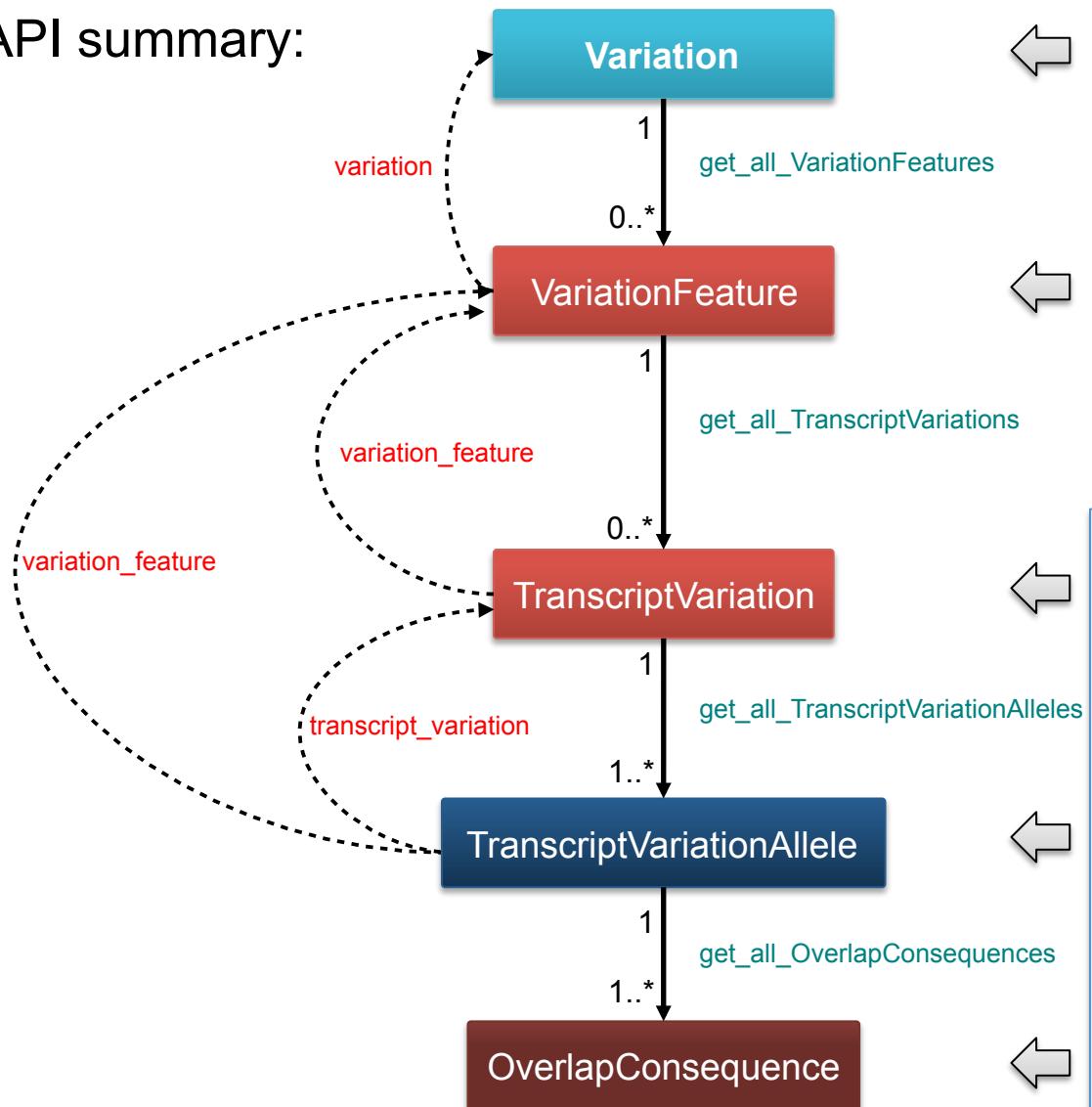
CAGCATCACTGCTGCAGCAACATCAGCAGCAGCATCAGCAACTCAGCATCACTGCTGCAGCAACATCAGCAGCAGCATCAGCA



→ http://www.ensembl.org/info/genome/variation/predicted_data.html#consequence_type_table

Transcript variation

API summary:



Variant information

Variant information at the genomic level

VariationFeature information at the transcript level

Specific allele in the transcript and its consequence(s)

Consequence(s) terms

Transcript variation object

- Transcript variation = an instance of a variation feature inside or near a transcript (i.e. can be more than one per variation feature → e.g. alternative splicing)
- Most “severe” consequence string is stored in the VariationFeature object:

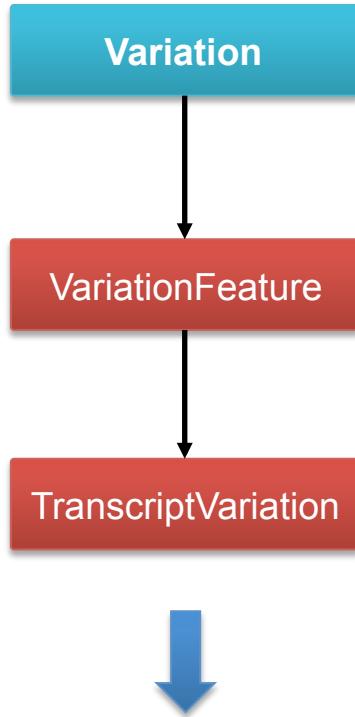
```
$vf->display_consequence()
```

- Can be retrieved from the **VariationFeature** objects as well as its object adaptor

Attribute	Example value(s)	Method(s)	Comment(s)
Transcript	Bio::EnsEMBL::Transcript	\$tv->transcript()	Returns a Core API object
Consequence type	intron_variant, stop_lost	\$tv->display_consequence() \$tv->consequence_type()	Only the most severe type All types in an array
Coordinates		\$tv->cdna_start() \$tv->cdna_end() \$tv->cds_start() \$tv->translation_start()	
Amino acid	F/L, */W	\$tv->pep_allele_string()	

Transcript variation

Example:



Name: [rs201036189](#)

Location: 1:206342907
 Allele string: C/T

Transcript: [ENST00000573034](#)
 Codon: Cgg/Tgg
 Amino acid string: R/W

Transcript (strand)	Consequence Type	Position in transcript	Position in CDS	Position in protein	Amino acid	Codons
ENST00000573034 (+) biotype: protein_coding	Missense variant	384	322	108	R/W	CGG/TGG

Transcript variation allele object

- Consequences are actually established at the allele level
- These represent the most specific information available

```
$transcript_variation->get_all_TranscriptVariationAlleles()
```

- All the consequences of a **TranscriptVariationAllele** are represented as a list of **OverlapConsequence** objects

Attribute	Example value(s)	Method(s)
TranscriptVariation	Bio::EnsEMBL::Variation::Transcri ptVariation	\$tva-> transcript_variation() *NB returns an object
OverlapConsequences	Reference list of Bio::EnsEMBL::Variation::OverlapC onsequence	\$tva-> get_all_OverlapConsequences() *NB returns a list of objects
HGVS notation at various levels	1:g.206516261C>T ENST00000295713.5:c.62C>T ENSP00000295713.5:p.Arg22Trp	\$tva-> hgvs_genomic() \$tva-> hgvs_coding() \$tva-> hgvs_protein()
Affected codons	Cgg/Tgg	\$tva-> display_codon_allele_string()
SIFT and PolyPhen predictions	string	\$tva-> sift_prediction() \$tva-> polyphen_prediction()

Protein function predictions

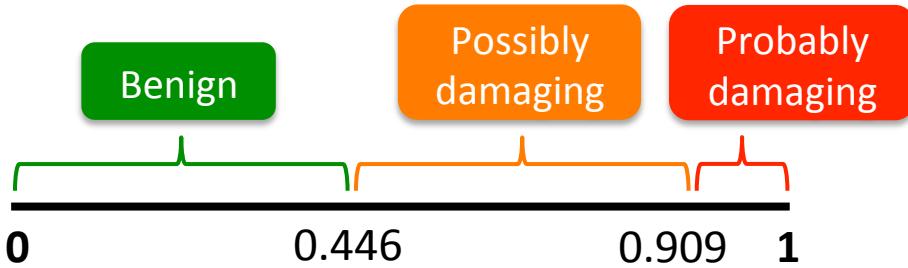
We run tools to identify non-synonymous mutations that are likely to affect protein function

SIFT

- Uses sequence homology and amino acid similarity to calculate if the substitution is: *tolerated* (>0.05) or *deleterious* (≤ 0.05)
- Supported species: Human, Mouse, Chicken, Cow, Dog, Horse, Pig, Rat, Sheep, Zebrafish

PolyPhen

- Uses sequence homology, PDB 3D structures, Pfam annotation, etc. to predict if a substitution is: *benign*, *possibly damaging*, *probably damaging* or *unknown*
- Supported species: human only

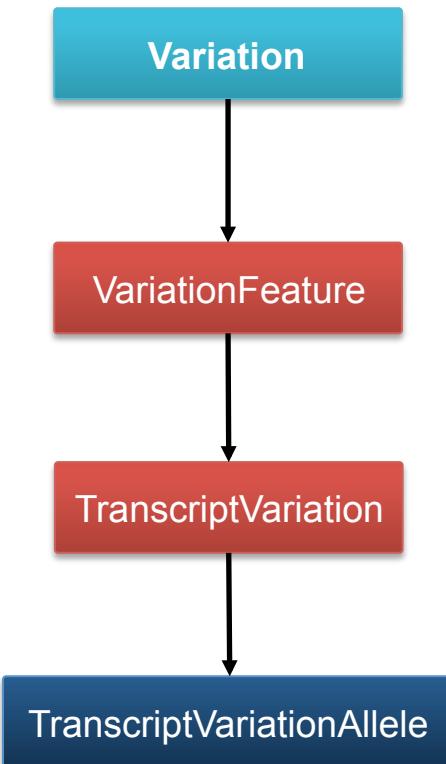


- Methods on a *TranscriptVariationAllele*:

- `$tva->sift_prediction()`
- `$tva->polyphen_prediction()`
- `$tva->sift_score()`
- `$tva->polyphen_score()`

Transcript variation allele

Example 1:
with 2 alleles



Name: [rs201036189](#)

Location: 1:206342907

Allele string: C/T

Transcript: [ENST00000573034](#)

Codon: Cgg/Tgg

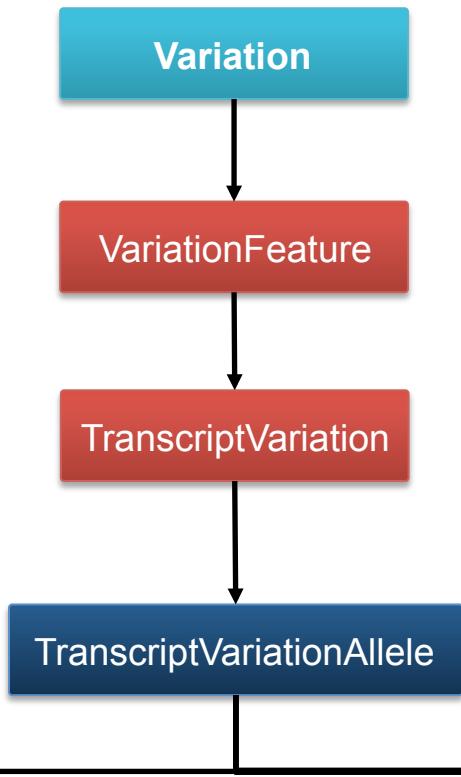
Amino acid string: R/W

Allele: T (alternative allele)

Transcript (strand)	Allele (transcript allele)	Consequence Type	Position in transcript	Position in CDS	Position in protein	Amino acid	Codons	SIFT	PolyPhen
ENST00000573034 (+) biotype: protein_coding	T (T)	Missense variant	384	322	108	R/W	CGG/TGG	0	0.998

Transcript variation allele (2)

Example 2:
with 3 alleles



Name: [rs191084422](#)

Location: 6:56553483
Allele string: C/G/T

Transcript: [ENST00000293925](#)
Codons: gCa/gGa/gTa
Amino acid string: A/G/V

Allele: G (alternative allele 1)
Allele: T (alternative allele 2)

Transcript (strand)	Allele (transcript allele)	Consequence Type	Position in transcript	Position in CDS	Position in protein	Amino acid	Codons	SIFT	PolyPhen
ENST00000293925 (+) biotype: protein_coding	G (G) 	Missense variant	1898	1898	633	A/G	GCA/GGA	0.01	0.117
ENST00000293925 (+) biotype: protein_coding	T (T) 	Missense variant	1898	1898	633	A/V	GCA/GTA	0.21	0.006

Overlap consequences

- Contain all the information we have about a specific consequence
- Also contains a ‘predicate’ which tests if this consequence applies to a **TranscriptVariationAllele**

Attribute	Example value(s)	Method(s)
Sequence Ontology term	missense_variant	\$oc-> so_term()
Label term	Missense variant	\$oc-> label()
Ensembl term	NON_SYNONYMOUS_CODING	\$oc-> display_term()
NCBI term	missense	\$oc-> NCBI_term()
Predicate	True/False	\$oc-> predicate(\$tva)

Exercise 8a

- Fetch all transcript variations in transcript [ENST00000001008](#) in human and retrieve the following:
 - variation name
 - consequence type (most severe)
 - amino acids *
 - position in cDNA *
 - position in translation *
- * if exist
- In a second attempt, filter for the transcript variations with the consequence type “missense_variant”

Hint: `fetch_all_by_Transcripts()` method requires a list reference of objects.

You have only one so use as parameter e.g.:

`fetch_all_by_Transcripts([$transcript])` instead of
`fetch_all_by_Transcripts($transcript)`

Exercise 8b

- Fetch all the coding transcript variation alleles in the transcript [ENST00000001008](#) in human and retrieve the following:
 - Allele string
 - Codon change (with the allele position displayed, e.g aAa)
 - Amino acid change
 - SIFT and PolyPhen predictions *

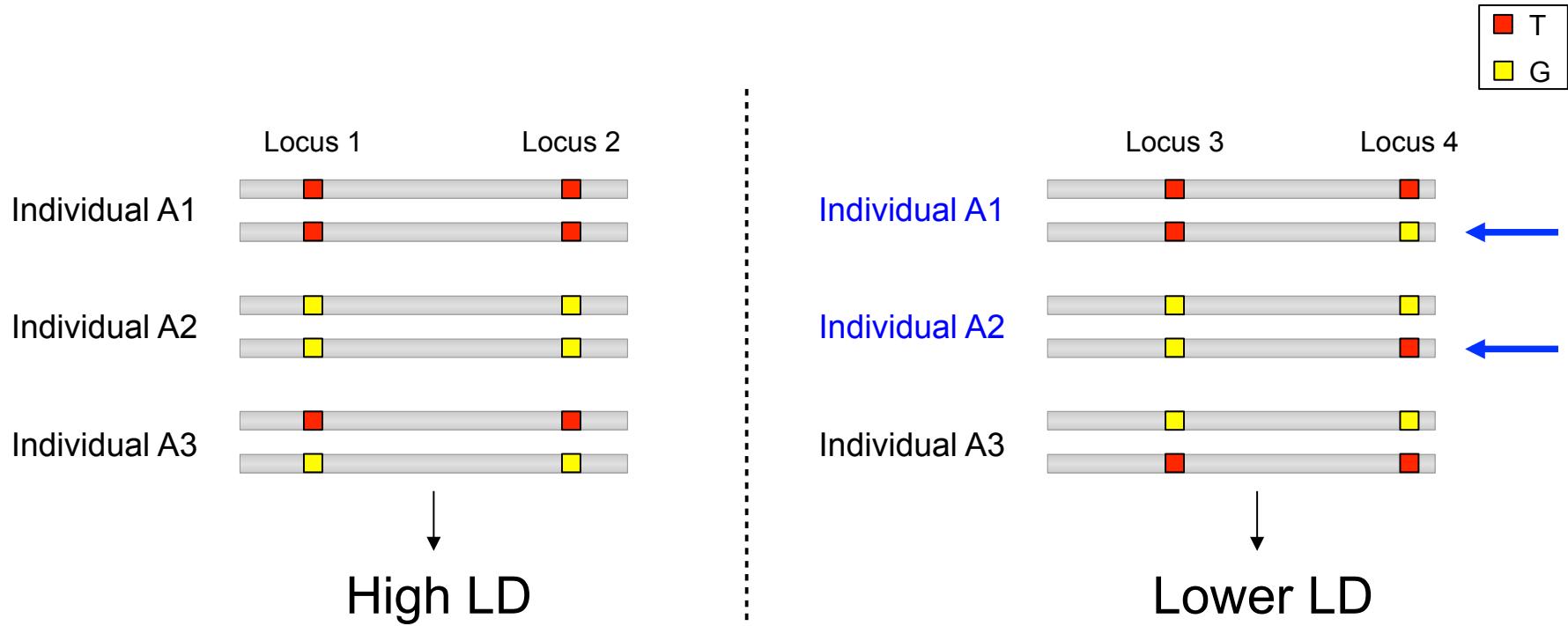
In this exercise, we only want information about the alternate allele.

Hint: You need to fetch the **TranscriptVariation** objects first, as you did in the exercise 8a

Linkage disequilibrium

Linkage disequilibrium (LD) is a measure of how frequently alleles at two separate loci are inherited together on the same haplotype in a specific population

- We provide LD statistics for the 1000 Genome Project phase 3 populations
- Two common measures:
 - r^2 , D' ($r^2 = 1 \rightarrow$ perfect LD)



Linkage disequilibrium container object



- Represents an instance of a container of pair-wise LD values in a region (2 by 2)
- Calculated on the fly (using a script written in C)
- Can contain values for multiple populations
- Most methods return hash references
e.g.: To retrieve the r^2 value you need to use:

```
$ldc->{ r2 }
```

Method
<code>\$ldc->get_all_ld_values ()</code>
<code>\$ldc->get_all_r_square_values ()</code>
<code>\$ldc->get_all_d_prime_values ()</code>

Key	Value
variation1	<i>VariationFeature</i> object of one of the two variations used for the calculation
variation2	<i>VariationFeature</i> object of the second variation used for the calculation
r2	0 to 1
d_prime	0 to 1

Exercise 9

Create an LD feature container for region chromosome **9:22124000-22126000** in human and find the names of all pairs of variants in high LD with each other in the 1000 Genomes ‘Toscani in Italy’ population, which is named “**1000GENOMES:phase_3:TSI**”

Hints:

- Try $r^2 \geq 0.95$
- The `fetch_by_slice()` method can take an optional argument “Population object”

Note:

Phase 3 Data from the 1000 Genomes Project is accessed directly from VCF files. It is not default API behaviour to read VCF’s, so to extract LD statistics for this population you will need to set the ‘`use_vcf`’ option on the database adaptor to 1.

```
$lda->db->use_vcf(1);
```



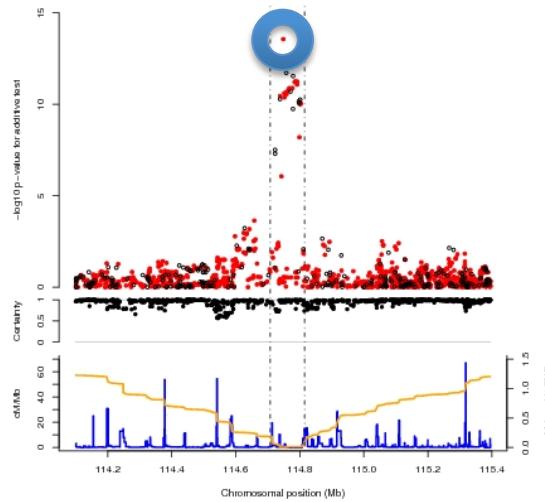
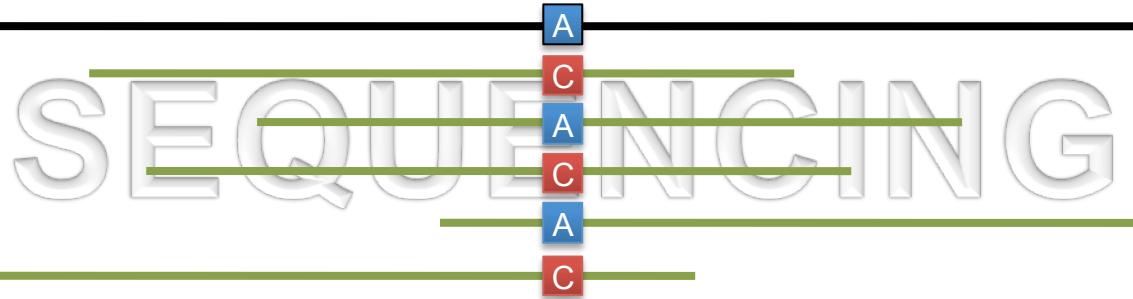
Variant *E*ffect *P*redictor

<http://www.ensembl.org/info/docs/tools/vep/index.html>

λ Background



Ref
Reads



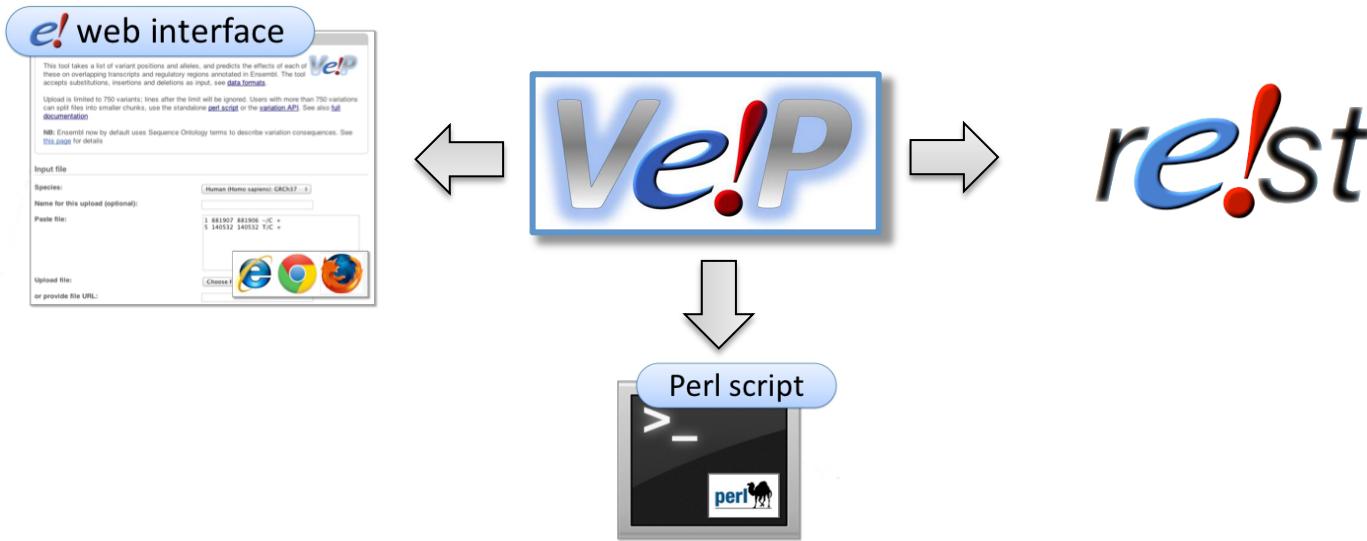
SNP



Variant Effect Predictor

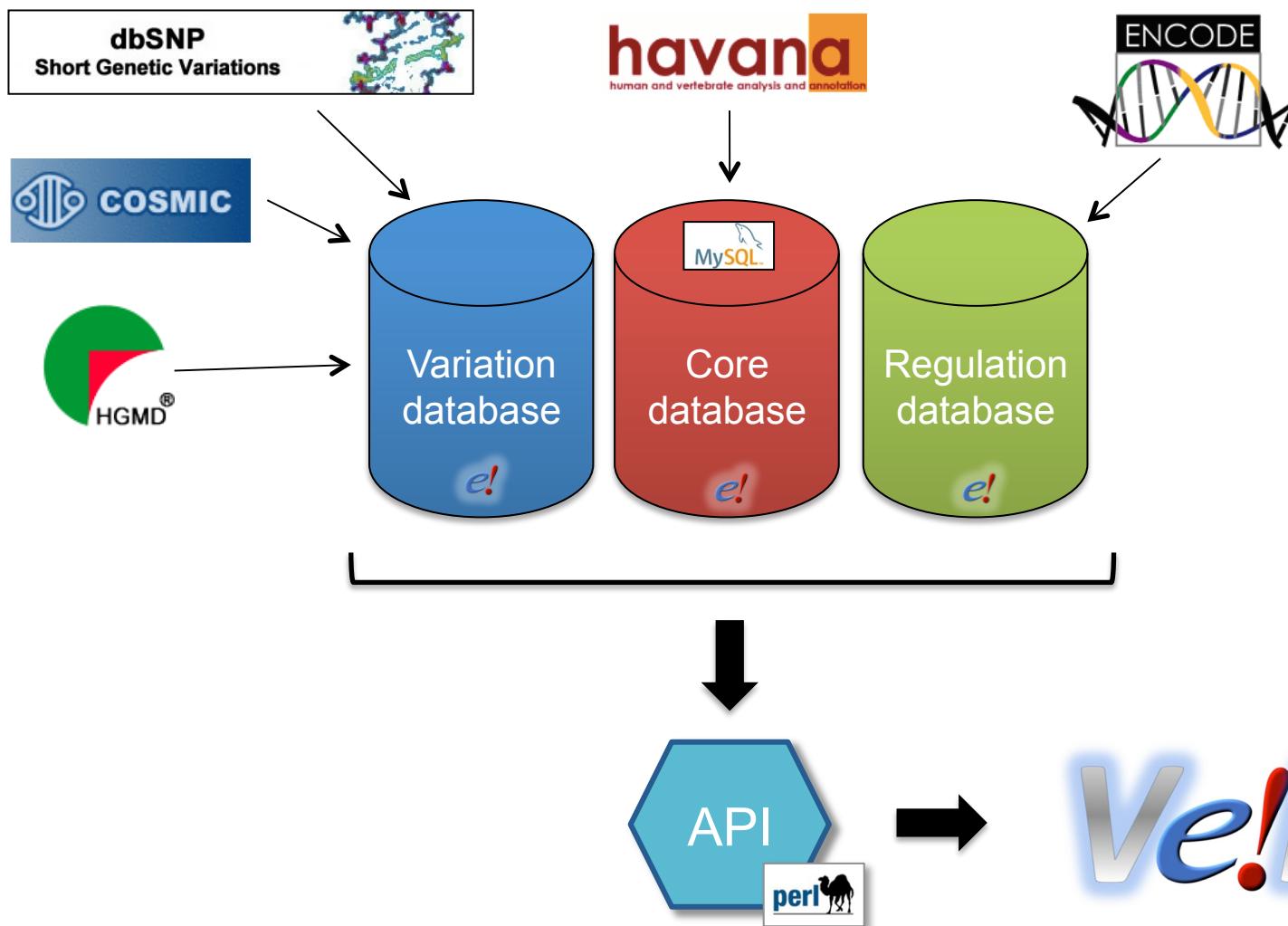


A tool to predict the functional consequences of variants, using the Ensembl API



- Use through simple web interface
- Use locally and completely privately by downloading the script and reference data
- Multiple input formats supported including VCF
- Many different annotations available
- Many filtering options available

Variant Effect Predictor



Variant Effect Predictor



- Uses all the consequences types from Ensembl Variation
- Can use aligned RefSeq transcripts (uses the Ensembl gene set by default)
- Provides variant descriptions using [HGVS](#) nomenclature
- Can use SIFT and PolyPhen2 predictions for popular species
- Provides regulatory region consequences
- Supports some types of Structural Variants encoded in VCF format
- Provides information on known variants at the same location:
 - Allele frequencies in reference populations
 - Pubmed identifiers if the variant is cited



The web interface

Variant Effect Predictor: Web



ensembl.org



BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Search: for

e.g. [BRCA2](#) or [rat 5:62797383-63627669](#) or [coronary heart disease](#)

Browse a Genome

The Ensembl project produces genome databases for vertebrates and other eukaryotic species, and makes this information freely available online.

Popular genomes



Human

GRCh38.p5



Human

GRCh37



Mouse

GRCm38.p4



Zebrafish

GRCz10

Still using Human GRCh37?



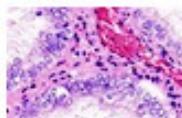
Go to



Variant Effect Predictor



Gene expression in different tissues



Find SNPs and other variants for my gene

GTRTATACTTC
CRTRAAAGTCTT
CTTCTAAATTCT
GRAACATTTGCC

Variant Effect Predictor: Web



ensembl.org/vep

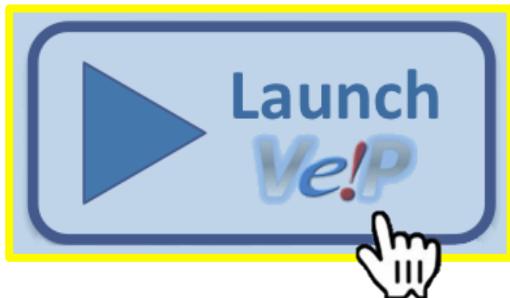
Help & Documentation > API & Software > Ensembl Tools > Variant Effect Predictor

Variant Effect Predictor



The VEP determines the effect of your variants (SNPs, insertions, deletions, CNVs or structural variants) on genes, transcripts, and protein sequence, as well as regulatory regions. Simply input the coordinates of your variants and the nucleotide changes to find out the:

- **genes and transcripts** affected by the variants
- **location** of the variants (e.g. upstream of a transcript, in coding sequence, in non-coding RNA, in regulatory regions)
- **consequence** of your variants on the protein sequence (e.g. stop gained, missense, stop lost, frameshift)
- **known variants** that match yours, and associated minor allele frequencies from the **1000 Genomes Project**
- SIFT and PolyPhen scores for changes to protein sequence
- ... And [more!](#)



Web interface

- Point-and-click interface
 - Suits smaller volumes of data
- [Documentation](#)
 [Launch the web interface](#)



Standalone perl script

- More options, more flexibility
 - For large volumes of data
- [Documentation](#)
 [Download latest version](#)



REST API

- Language-independent API
 - Simple URL-based queries
 - GET single variants, POST many
- [Documentation](#)

Variant Effect Predictor: Web



Variant Effect Predictor ?

Click for help

i VEP for Human GRCh37

If you are looking for VEP for Human GRCh37, please go to [GRCh37 website](#).

Species:

Human (Homo sapiens)



Assembly: GRCh38.p5

Select species

Name for this data (optional):

Either paste data:

Input data

Examples: [Ensembl default](#), [VCF](#), [Variant identifiers](#), [HGVS notations](#), [Pileup](#)

Choose File No file chosen

Upload data
• 50MB limit
• compressed OK

Or upload file:

Or provide file URL:

Transcript database to use:

Ensembl transcripts

Gencode basic transcripts

RefSeq transcripts

Ensembl and RefSeq transcripts

Include additional EST and CCDS transcripts:



Identifiers and frequency data Additional identifiers for genes, transcripts and variants; frequency data

Extra options e.g. SIFT, PolyPhen and regulatory data

Filtering options Pre-filter results by frequency or consequence type

Run >

Clear

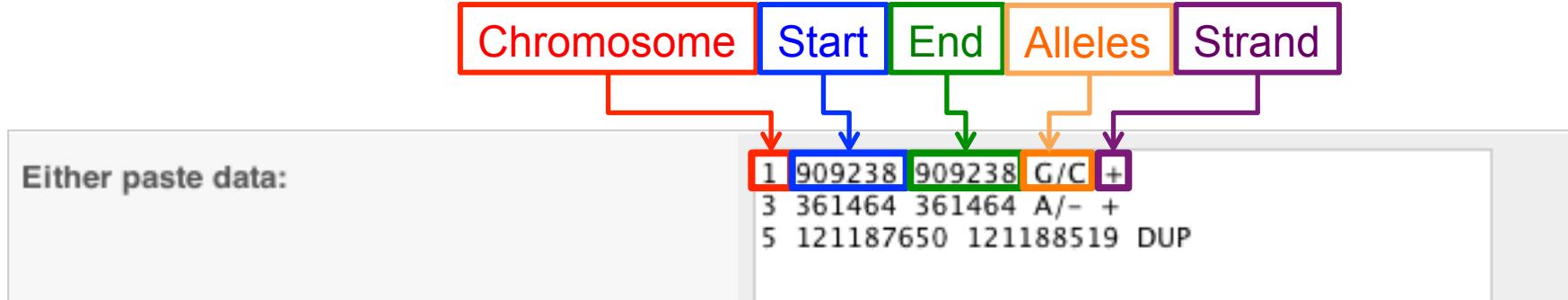
Submit job

Variant Effect Predictor: Web



Input format

ensembl.org/info/docs/tools/vep/vep_formats.html#input



- Coordinates
 - 1-based
 - for a SNP start = end
- Alleles
 - “/” separated list (may be more than two), reference allele first!
 - “-” represents no sequence (insertion, deletion)
- Other formats supported:
 - VCF
 - HGVS nomenclature
 - Variant identifiers (e.g. rsIDs from dbSNP)

Variant Effect Predictor: Web



Jobs

ensembl.org/info/docs/tools/vep/online/input.html#jobs

Refresh

Analysis	Ticket	Jobs	Submitted at	Filter
Variant Effect Predictor	XNMrwWEkBBxhoSg	Job 1: VEP analysis of pasted data in Homo_sapiens Queued Done View Results	10/03/2014, 16:57	
Variant Effect Predictor	KsInwPcsbOlfiZAr	Job 1: VEP analysis of pasted data in Homo_sapiens	10/03/2014, 16:57	

Ticket identifier

Job name

- Queued** - the job is waiting in line to run
- Running** - the job is running
- Done** - the job is finished and you can view the results
- Failed** - the job has failed for some reason

Save to your account (if logged in)

Edit and resubmit your job

Delete job

Variant Effect Predictor: Web



Return to jobs



BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Search all species

Human (GRCh37) ▾

Location: 6:133,002,729-133,035,188 ▾

Gene: VNN1 ▾

Transcript: VNN1-001 ▾

Variation: rs143624951 ▾

Jobs

Variation displays

- Explore this variation
- Genomic context
 - Genes and regulation (1)
 - Flanking sequence
 - Population genetics
 - Individual genotypes
 - Linkage disequilibrium
 - Phenotype Data
 - Phylogenetic Context (6)
 - Citations
- External Data
 - SNPedia
 - LOVD

rs143624951 SNP

Original source

Variants (including SNPs and indels) imported from dbSNP (release 138) | [View in dbSNP](#)

Alleles

C/T | Ancestral: C | Ambiguity code: Y

Location

Chromosome 6:133004370 (forward strand) | [View in location tab](#)

Most severe consequence

| Stop gained | [See all predicted consequences \[Genes and regulation\]](#)



Evidence status ⓘ

HGVS names +

This variation has 3 HGVS names - click the plus to show

Genes and regulation ⓘ

Gene and Transcript consequences

Show/hide columns

Gene	Transcript (strand)	Allele (transcript allele)	Consequence Type	Position in transcript	Position in CDS	Position in protein	Amino acid	Code
ENSG00000112299	ENST00000367928	T (A)	Stop gained	1465	1451	484	W/*	TGG
HGNC: VNN1	(-)	biotype: protein_coding						



Variant Effect Predictor: Web



Results summary

 ensembl.org/info/docs/tools/vep/online/results.html#summary

Variant Effect Predictor results

Summary statistics for ticket uTPZdaCscZ7TeOIX: □

Category	Count
Variants processed	19976
Variants remaining after filtering	19976
Novel / existing variants	914 (4.6%) / 19062 (95.4%)
Overlapped genes	116
Overlapped transcripts	428
Overlapped regulatory features	373

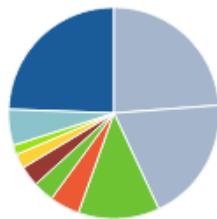
Variant counts

- before/after filtering
- novel / known variants

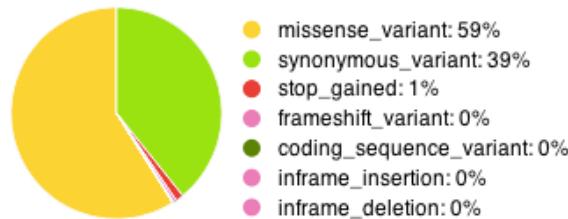
Consequence type charts

- total observed - may be more than one per variant
- Separate chart for coding consequences

Consequences (all)



Coding consequences



Variant Effect Predictor: Web



Results preview



ensembl.org/info/docs/tools/vep/online/results.html#table

Navigate results

NB one row per variant/
transcript overlap!

Create and edit filters

NB more columns:
scroll right!

Results preview

Navigation

Page: < < 1 of 1 > > | Show: All variants

Filters

Uploaded variant is defined

Download

All: VCF VEP TXT

BioMart: Variants Genes

Show/hide columns

Uploaded variant	Location	Allele	Consequence	Impact	Symbol	Gene	Feature type	Feature ID	Score	Effect
rs116383664	1:118008 1180081		upstream_gene_variant	MODIFIER	TTLL10-AS1	ENSG00000205231	Transcript	ENST00000205231	379317	antisense
rs116383664	1:1180081 1180081		upstream_gene_variant	MODIFIER	TTLL10	ENSG00000162571	Transcript	ENST00000162571	36379	nonsense-mediated_
rs116383664	1:1180081- 1180081		missense_variant	MODERATE	TTLL10	ENSG00000162571	Transcript	ENST00000162571	290	protein_coding

Scroll to see more
columns »

Show/hide columns
in results table

- Download results
- Send results to BioMart

VEP - Exercise



- Using a nucleotide change of G → A at position 25587758 on chromosome 21:
 - a) Does this variant cause an amino acid change?
 - b) Is there an existing variant here?
 - c) What is its global minor allele frequency (GMAF)?
 - d) Which do you think is the “worst” consequence listed?
 - e) What is the HGNC symbol of the gene affected?



The script version



Ressources

- Documentation

<http://www.ensembl.org/info/docs/tools/vep/script/index.html>

- Download

http://www.ensembl.org/info/docs/tools/vep/script/vep_download.html

- Help

dev@ensembl.org (public - mailing list)

helpdesk@ensembl.org (private - ticketed helpdesk)



Advantages compared to the web version

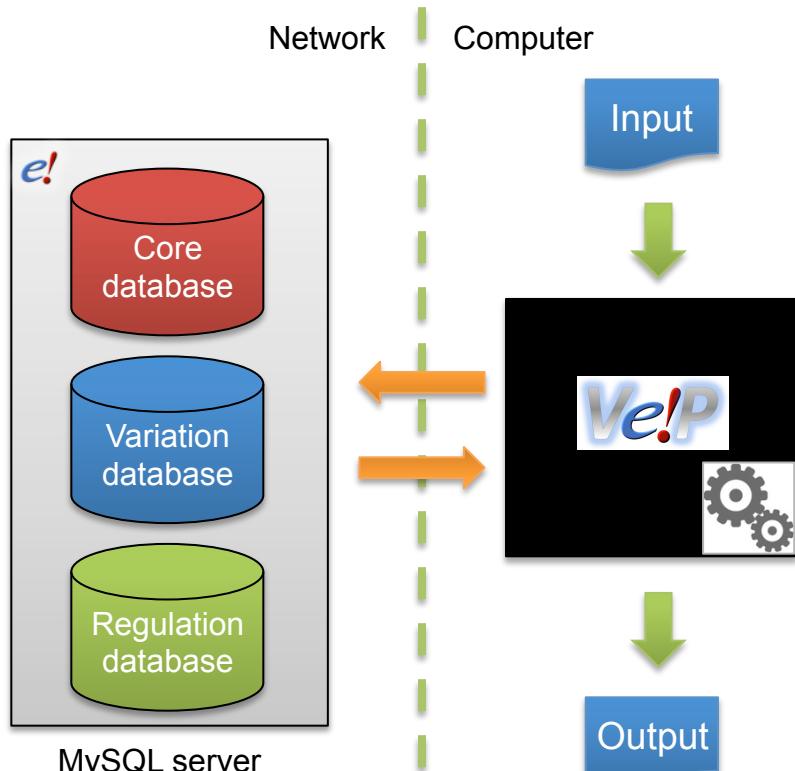
- Can be included in a pipeline
- No input file limit
- More advanced use (i.e. more filters and options)
- The possibility to use existing plugins or to write your own plugins
You can see the list of available plugins here:
https://github.com/Ensembl/VEP_plugins
- The possibility to run the script on stand alone, without internet connection, using the VEP cache files:
http://www.ensembl.org/info/docs/tools/vep/script/vep_cache.html

Variant Effect Predictor: Script

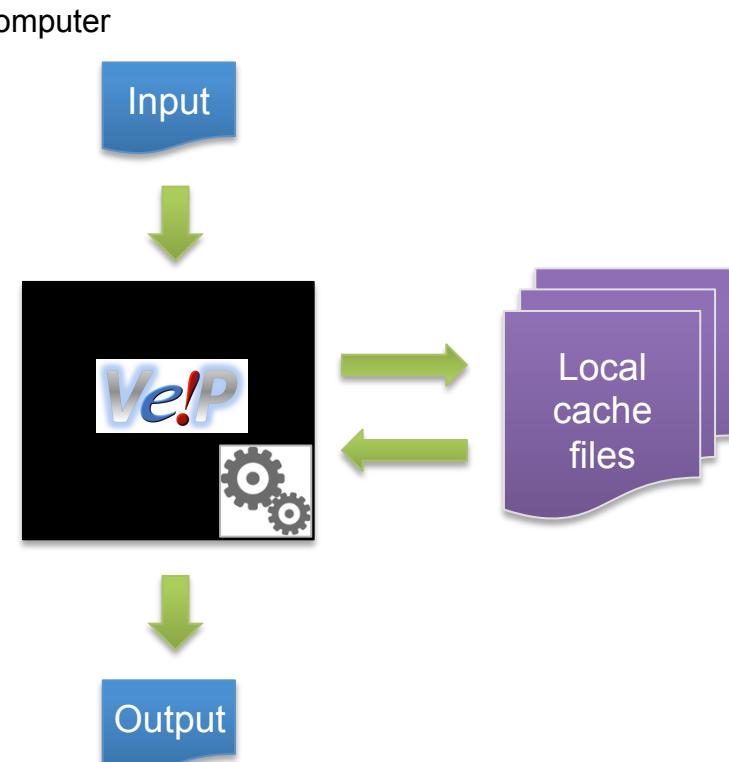


Online vs offline

Online “normal” mode



Offline “cache” mode



Getting more information

- Variation documentation

<http://www.ensembl.org/info/genome/variation/index.html>

<http://www.ensembl.org/info/docs/api/variation/index.html>

- Online Perl API documentation

<http://www.ensembl.org/info/docs/Doxygen/variation-api/index.html>

- Variation API tutorial

http://www.ensembl.org/info/docs/api/variation/variation_tutorial.html

- Ensembl developers mailing list

dev@ensembl.org

Acknowledgements

The Entire Ensembl Team

Andrew Yates¹, Wasiu Akanni¹, M. Ridwan Amode¹, Daniel Barrell^{1,2}, Konstantinos Billis¹, Denise Carvalho-Silva¹, Carla Cummins¹, Peter Clapham², Stephen Fitzgerald¹, Laurent Gil¹, Carlos García Girón¹, Leo Gordon¹, Thibaut Hourlier¹, Sarah E. Hunt¹, Sophie H. Janacek¹, Nathan Johnson¹, Thomas Juettemann¹, Stephen Keenan¹, Ilias Lavidas¹, Fergal J. Martin¹, Thomas Maurel¹, William McLaren¹, Daniel N. Murphy¹, Rishi Nag¹, Michael Nuhn¹, Anne Parker¹, Mateus Patricio¹, Miguel Pignatelli¹, Matthew Rahtz², Harpreet Singh Riat¹, Daniel Sheppard¹, Kieron Taylor¹, Anja Thormann¹, Alessandro Vullo¹, Steven P. Wilder¹, Amonida Zadissa¹, Ewan Birney¹, Jennifer Harrow², Matthieu Muffato¹, Emily Perry¹, Magali Ruffier¹, Giulietta Spudich¹, Stephen J. Trevanion¹, Fiona Cunningham¹, Bronwen L. Aken¹, Daniel R. Zerbino¹ and Paul Flicek^{1,2,*}

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK and ²Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

Funding



National
Human Genome
Research Institute



Centre for Therapeutic
Target Validation



epigenome



Co-funded by the
European Union



EMBL-EBI



Feedback

https://www.surveymonkey.co.uk/r/API_Jan16

Figures

- Babies: <http://www.impawards.com/2010/posters/babies.jpg> (slide #2)