



IX International Course of Massive Data Analysis FOR GENOMICS



Where are we?



Sequence preprocessing

Mapping

Variant Calling

Variant prioritization

Functional annotation

Gene-Set Analysis

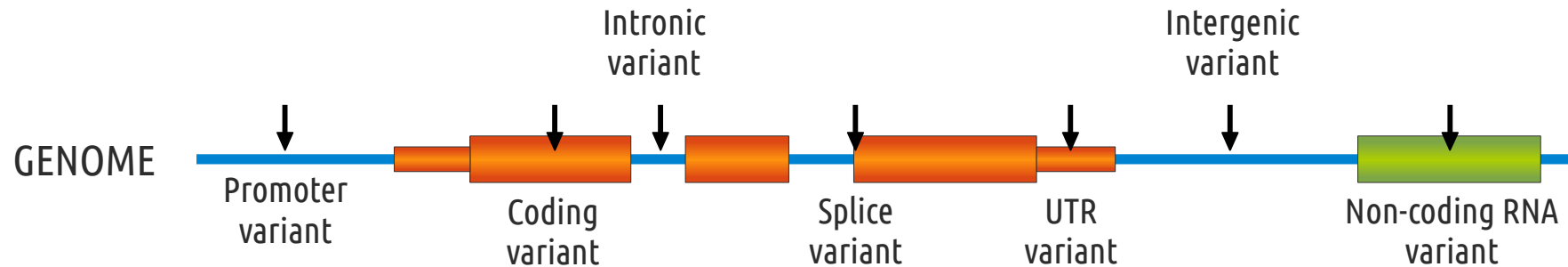
GWAS Analysis



The haystack has been cleared away to reveal a large pile of needles



What is functional annotation?



Why we do that?

- ▶ Each individual exome carries ~25,000 variants → **PRIORITIZATION!**
- ▶ We want to identify a **small subset** of functionally important variants to pinpoint the putative disease causal variants
- ▶ We need strategies to **estimate the deleteriousness** of our variants to better identify disease-causal variants

CAUTION!

On average, each *normal* person is found to carry:

~11,000 synonymous variants

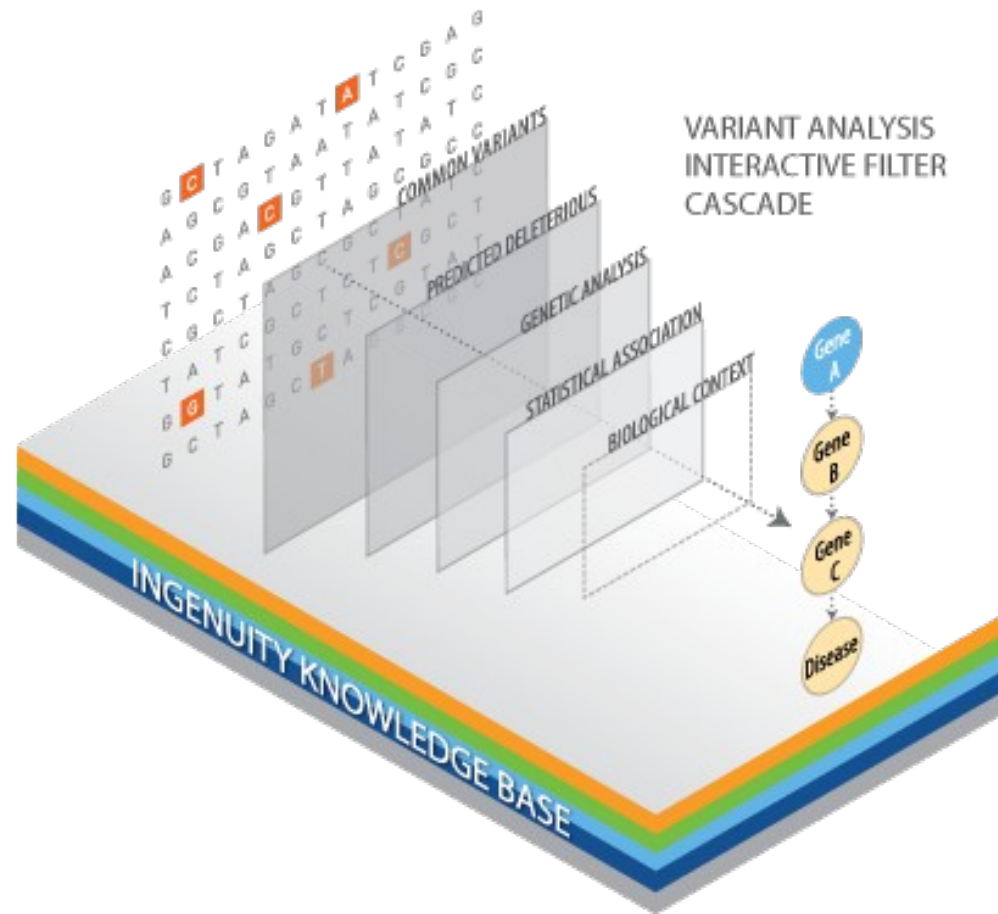
~11,000 non-synonymous variants

250 to 300 **los-of-function** variants in annotated genes

50 to 100 variants previously implicated in **inherited disorders**

1000 Genomes Project Consortium. *A map of human genome variation from population-scale sequencing.* **Nature.** 2010 Oct 28;467(7319):1061-73. PubMed PMID: 20981092

The objective



Sources of functional information

Table 1 Publicly available tools and databases for various tasks of genetic variant annotation and prioritization

Category	Database/tool/project	Description	URL
Genetic variant data sources	dbSNP ⁶⁸	Comprehensive, curated SNP and short indel database	http://www.ncbi.nlm.nih.gov/projects/SNP
	DbVar ⁶⁹	Comprehensive, curated database for structural variants	http://www.ncbi.nlm.nih.gov/dbvar
	DGV ⁷⁰	Human structural variants from samples with no phenotype	http://projects.tcag.ca/variation
Functional characterization of genomic elements	ENCODE ⁷¹	High-throughput functional characterization of DNA elements, including noncoding regions	http://www.genome.gov/10005107
	SIFT ⁷² , PolyPhen ⁷³	Prioritization of nonsynonymous SNPs	http://sift.jcvi.org , http://genetics.bwh.harvard.edu/pph2
Public gene–trait associations	dbGaP ³⁴	Comprehensive listing of genotype-to-phenotype mappings	http://www.ncbi.nlm.nih.gov/gap
	EGA ⁷⁴	Genotype–phenotype experiment archive	http://www.ebi.ac.uk/ega
Disease-associated mutations	HGMD ³⁵	Database for human disease mutations	http://www.hgmd.org
	OMIM ³⁶	Mendelian disease gene associations	http://www.ncbi.nlm.nih.gov/omim
	SwissVar ⁷⁶	Variant catalog of the UniProt knowledge bases	http://swissvar.expasy.org
	GAD ⁷⁷	NCBI source for genotype–disease associations	http://geneticassociationdb.nih.gov
	GWAS catalog from NHGRI ⁷⁸	SNP-phenotype associations found by GWAS	http://www.genome.gov/gwastudies
Whole-genome repositories	Complete genomics public genomes ⁷⁹	Complete genomics for 69 genomes from multiple ancestries (includes samples from the NHGRI and NIGMS repositories)	http://www.completegenomics.com/sequence-data/download-data
	1,000 Genomes ⁸⁰	Expanding resource currently housing three low-coverage whole genomes of multiple ancestries	http://www.1000genomes.org
Ancestry-focused variant data sources	HapMap ²⁶	Haplo-block mapping for diverse populations	http://www.hapmap.org
	HGDP ²⁷	SNP profiles of samples from several endogenous populations	http://hagsc.org/hgdp
Pharmacogenomic associations and data sources	PharmGKB ⁵⁶	Variant–pharmacokinetic/pharmacodynamic trait associations and gene–drug interactions	http://www.pharmgkb.org
	DrugBank ⁸¹	Drug-target database with biochemical properties	http://drugbank.ca



Cordero P, Ashley EA. *Whole-genome sequencing in personalized therapeutics*. *Clin Pharmacol Ther*. 2012 Jun ;91(6):1001-9. PubMed PMID: 22549284

Sources of functional information - ENCODE

ENCODE (The Encyclopedia of DNA Elements)

- ▶ Objective: Characterize all of the functional elements encoded in the human genome.
- ▶ The 80% of the human genome has a defined biochemical function in at least one cell type.
- ▶ Focus on non-coding regions: accessible DNA regions, DNA-binding motifs and RNA-coding.

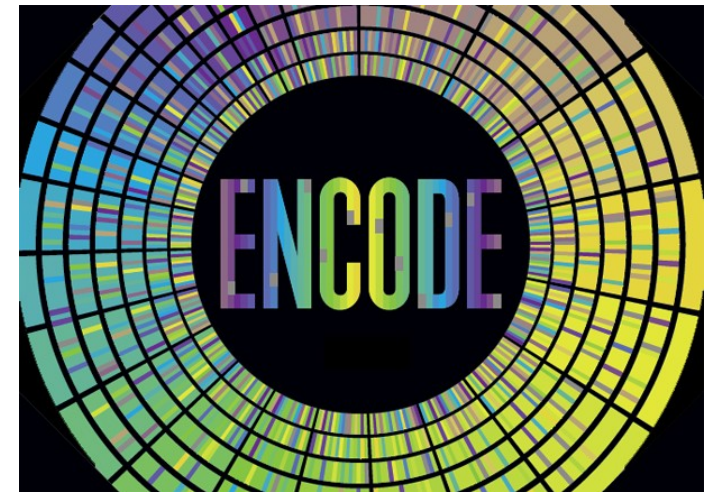
ARTICLE

doi:10.1038/nature11247

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.



The ENCODE Project Consortium. *An integrated encyclopedia of DNA elements in the human genome.* *Nature*. 2012 Sep 6;489(7414):57-74. Pubmed PMID: 22955616

Computational methods and tools

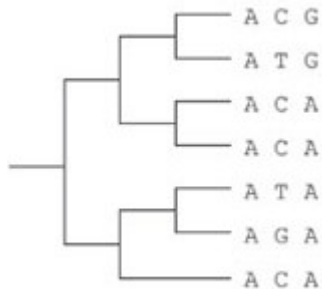
- ▶ **Annotated information** is sometimes **limited**, particularly for rare and complex traits
- ▶ Computational methods can measure deleteriousness by using **comparative genomics** and knowledge of **protein biochemistry and structure**

Comparative Genomics

Focus on sequences that have not been removed by **natural selection**.

Quantify evolutionary changes in genes or genomes and define conserved and neutral regions.

Variants observed in conserved sites are highly likely to be **deleterious**.



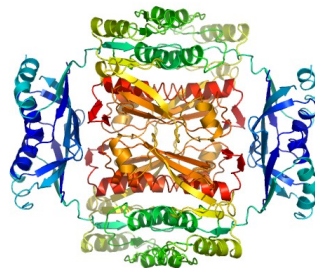
Effects in protein-coding variants

Can combine **evolutionary** and **biochemical** information.

Two types of approaches: first **principles** approaches and **trained** approaches.

Use **alignments of homologous proteins** to estimate mutational deleteriousness.

Use **biochemical data** such as amino acid properties, binding information and structural information to estimate the impact.



Effects in non-coding variants

The majority of the human genetic variation is in non-coding regions.

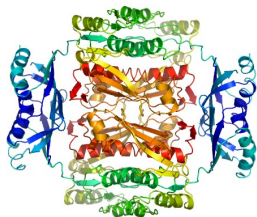
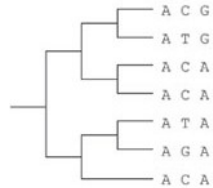
No detectable conservation outside vertebrates.

Main strategy for estimation is testing the **mammalian conservation** of the non-coding variants.

Cooper GM, Shendure J. *Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data.* **Nature Reviews Genetics.** 2011 Aug 18;12(9):628-40. Pubmed PMID: 21850043

Computational methods and tools

Prediction scores for non-synonymous variants



Cooper GM, Shendure J. *Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data.* **Nature Reviews Genetics.** 2011 Aug 18;12(9):628-40. Pubmed PMID: 21850043

Table 1 | **Tools for protein-sequence-based prediction of deleteriousness**

Name	Type	Information	URL	Refs
MAPP	Constraint-based predictor	Evolutionary and biochemical	http://mendel.stanford.edu/SidowLab/downloads/MAPP/index.html	27
SIFT	Constraint-based predictor	Evolutionary and biochemical (indirect)	http://sift.bii.a-star.edu.sg/	39
PANTHER	Constraint-based predictor	Evolutionary and biochemical (indirect)	http://www.pantherdb.org/	41
MutationTaster*	Trained classifier	Evolutionary, biochemical and structural	http://www.mutationtaster.org/	40
nsSNP Analyzer	Trained classifier	Evolutionary, biochemical and structural	http://snpanalyzer.uthsc.edu/	44
PMUT	Trained classifier	Evolutionary, biochemical and structural	http://mmb2.pcb.ub.es:8080/PMut/	38
polyPhen	Trained classifier	Evolutionary, biochemical and structural	http://genetics.bwh.harvard.edu/pph2/	35
SAPRED	Trained classifier	Evolutionary, biochemical and structural	http://sapred.cbi.pku.edu.cn/	42
SNAP	Trained classifier	Evolutionary, biochemical and structural	http://www.rostlab.org/services/SNAP/	36
SNPs3D	Trained classifier	Evolutionary, biochemical and structural	http://www.snps3d.org/	51
PhD-SNP	Trained classifier	Evolutionary and biochemical (indirect)	http://gpcr2.biocomp.unibo.it/~emidio/PhD-SNP/PhD-SNP_Help.html	37

*Also makes predictions for synonymous and non-coding variant effects: for example, splicing. MAPP, Multivariate Analysis of Protein Polymorphism; polyPhen, polymorphism phenotyping.

Computational methods

Prediction scores for non-coding variation

Table 2 | **Tools for nucleotide-sequence-based prediction of deleteriousness**

Name	Type	Information	URL	Refs
phastCons	Phylogenetic HMM	Evolutionary	http://compgen.bscb.cornell.edu/phast/	60
GERP	Single-site scoring	Evolutionary	http://mendel.stanford.edu/SidowLab/downloads/gerp/index.html	67
Gumby	Single-site scoring	Evolutionary	http://pga.jgi-psf.org/gumby/	21
phyloP	Single-site scoring	Evolutionary	http://compgen.bscb.cornell.edu/phast/	66
SCONE	Single-site scoring	Evolutionary	http://genetics.bwh.harvard.edu/scone/	68
binCons	Sliding-window scoring	Evolutionary	http://zoo.nhgri.nih.gov/binCons/index.cgi	69
Chai Cons	Sliding-window scoring	Evolutionary and structural	http://research.nhgri.nih.gov/software/chai	71
VISTA	Visualization tool (various scores)	Evolutionary	http://genome.lbl.gov/vista/index.shtml	70

GERP, Genomic Evolutionary Rate Profiling; HMM, hidden Markov model; SCONE, Sequence Conservation Evaluation.

Cooper GM, Shendure J. *Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data.* **Nature Reviews Genetics.** 2011 Aug 18;12(9):628-40. Pubmed PMID: 21850043

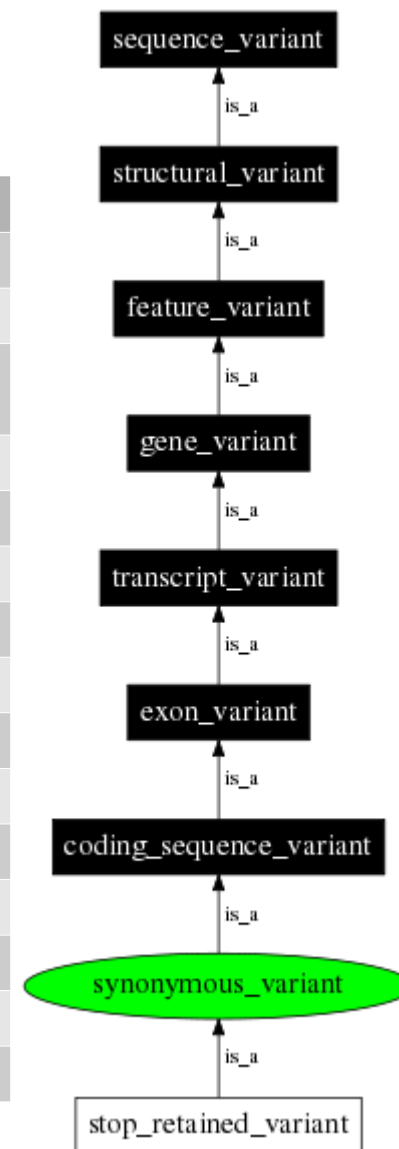
Consequence types

The standard:



Common vocabulary

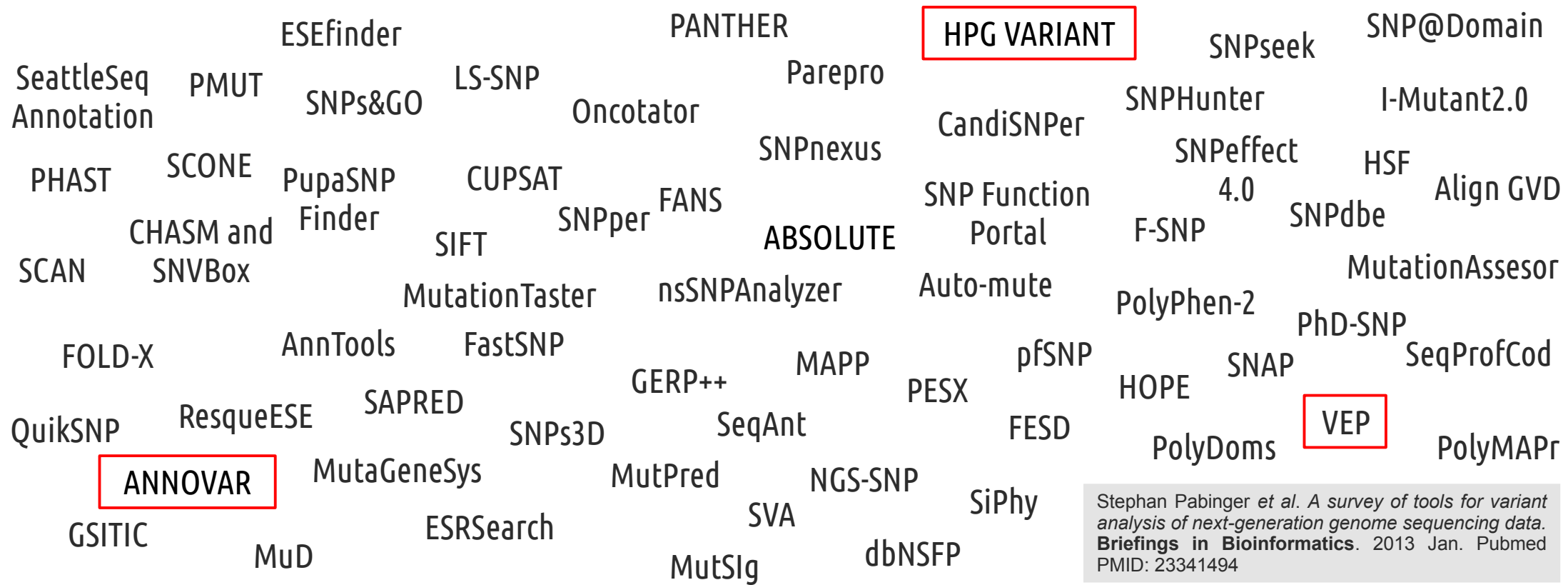
Label	SO accession	Description
Coding sequence	SO:0001580	In coding sequence with in determinate effect
Synonymous codon	SO:0001588	In coding sequence, not resulting in an amino acid change (silent mutation)
Non-synonymous codon	SO:0001583	In coding sequence and results in an amino acid change in the encoded peptide sequence
Stop gained	SO:0001587	In coding sequence, resulting in the gain of a stop codon
Stop lost	SO:0001587	In coding sequence, resulting in the gain of a stop codon
Splice site	SO:0001630	1-3bps in to an exon or 3-8bps into an intron
Splice acceptor	SO:0001574	A splice variant that changes the 2 base region at the 3' end of an intron
Splice donor	SO:0001575	A splice variant that changes the 2 base region at the 5' end of an intron
5' UTR	SO:0001623	In 5 prime untranslated region
3' UTR	SO:0001624	In 3 prime untranslated region
Upstream	SO:0001635	Within 5kb upstream of the 5 prime end of a transcript
Downstream	SO:0001633	Within 5kb downstream of the 3 prime end of a transcript
TFBS	SO:0001782	A sequence variant located with in a transcription factor binding site
miRNA target	SO:0000934	A binding site where the molecule is a microRNA
Intergenic	SO:0001628	More than 5 kb either upstream or downstream of a transcript



More information: http://www.ensembl.org/info/docs/variation/predicted_data.html

Tools for functional annotation

- ▶ We need to measure the **impact** of each variant in the genome
- ▶ We **cannot** annotate 25,000 variants **manually** checking more than 20 databases
- ▶ **Tools integrate** biological information and **ease** the functional annotation of hundreds of thousand variants



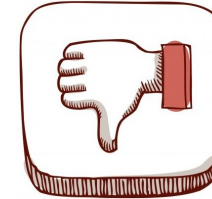
Stephan Pabinger *et al.* A survey of tools for variant analysis of next-generation genome sequencing data. *Briefings in Bioinformatics*. 2013 Jan. Pubmed PMID: 23341494

ANNOVAR

ANNOVAR web site: <http://www.openbioinformatics.org/annovar/>



- Free and open source
- Can annotate SNV, insertions and deletions
- **Regulatory information:** Conserved genomic regions, TFBSs, miRNA targets and predicted miRNA secondary structures. ENCODE DNase I hypersensitive sites, Histone methylations, CHIP and RNA-Seq peaks
- DbSNP, 1000 genomes SIFT and GERP filtering
- **Predictions:** Polyphen, LRT, MutationTaster, PhyloP
- Can handle **custom annotations** in GFF3
- Can handle 1 o 0-based coordinates
- **5 Species** (human, mouse, worm, fly, yeast)



- Accepts VCF4, GFF3-SOLiD and CSV BUT after conversion to their **particular input file:**

Chr	Start	End	Ref	Obs	Comments
1	161003	161003	C	T	comments: rs1000050

- **Perl** written program
- **Installation** required
- Users need to **download** every annotation database and save them locally (~35GB per assembly)
- Need to be **run several times**
- **Output:** several files depending on the query
- Does not use **Sequence Ontology** terms

Wang K, Li M, Hakonarson H. ANNOVAR: Functional annotation of genetic variants from next-generation sequencing data. **Nucleic Acids Research**. Sep;38(16):e164 Pubmed PMID: 20601685

ANNOVAR

EXAMPLE of ANNOVAR usage

DOWNLOADING BIOLOGICAL DATA:

```
user@computer:~$ annotate_variation.pl -buildver hg19 -downdb refgene humandb/
```

```
user@computer:~$ annotate_variation.pl -buildver hg19 -downdb snp135 -webfrom annovar humandb/
```

```
user@computer:~$ annotate_variation.pl -buildver hg19 -downdb phastConsElements46way humandb/
```

```
user@computer:~$ annotate_variation.pl -buildver hg19 -downdb 1000g2012apr -webfrom annovar  
humandb/
```

```
user@computer:~$ annotate_variation.pl -buildver hg19 -downdb cytoBand humandb/
```

EXTRACTING THE EFFECT:

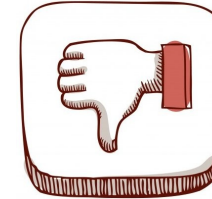
```
user@computer:~$ annotate_variation.pl -geneanno example/ex1.human humandb/
```

```
user@computer:~$ annotate_variation.pl -regionanno -dbtype band example/ex1.human humandb/
```

```
user@computer:~$ annotate_variation.pl -filter -dbtype 1000g2012apr_eur example/ex1.human  
humandb/
```


Variant Effect Predictor (VEP)

VEP documentation site: <http://www.ensembl.org/info/docs/variation/vep/index.html>



- Backed by **Ensembl**
- Free and open source
- **3 ways of functionality**: web interface, standalone Perl script and Ensembl's Perl API
- **Input** formats: CSV, VCF, Pileup and HGVS
- **Regulatory information**: TFBSs
- **Filtering** by coding regions and MAF
- **Predictions**: SIFP, PolyPhen
- 1000 genomes and dbSNP information
- Uses **Sequence Ontology**
- **Many species**
- Regulatory information does **not** include miRNA targets
- Web interface limit: **750 variants**
- The **standalone Perl script** needs:
 - **Perl** and **MySQL** support (more than 100GB of data)
 - **Download, install and update** every ~ 2 months
- Perl **API** requires:
 - **Installation** (Really, really hard!)
 - **Downloads and update**
 - API documentation → **Hard to understand**

McLaren W, Pritchard B, Rios D, Chen Y, Flicek P, Cunningham F. *Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor*. **BMC Bioinformatics** 26(16):2069-70(2010) Pubmed PMID: 20562413

Variant Effect Predictor (VEP)

VEP web interface: <http://www.ensembl.org/tools.html>

The screenshot shows the VEP web interface. On the left is a navigation menu with options like 'Add your data', 'Attach DAS', 'Manage Data', and 'Features on Karyotype'. The main area is titled 'Variant Effect Predictor:' and contains instructions: 'This tool takes a list of variant positions and alleles, and predicts the effects of each of these on overlapping transcripts and regulatory regions annotated in Ensembl. The tool accepts substitutions, insertions and deletions as input, see [data formats](#).' It also states: 'Upload is limited to **750 variants: lines after the limit will be ignored**. Users with more than 750 variations can split files into smaller chunks, use the standalone [perl script](#) or the [variation API](#). See also [full documentation](#).' A note says: 'NB: Ensembl now by default uses Sequence Ontology terms to describe variation consequences. See [this page](#) for details'. Below this is the 'Input file' section with fields for 'Species' (set to 'Human (Homo sapiens): GRCh'), 'Name for this data (optional)', 'Paste data:' (containing a table of variants), 'Upload file:' (with a 'Choose File' button), 'or provide file URL:', 'Input file format:' (set to 'Ensembl default'), and 'Options' (with radio buttons for 'Ensembl transcripts' and 'RefSeq and other transcripts').

Line	Chromosome	Start	End	Ref	Alt
1	881907	881906	-/C	+	
5	140532	140532	T/C	+	

The web interface to the VEP has a hard **limit of 750 variants** in your uploaded file. However, it is possible that the tool will not work with fewer variants than this, depending on the content of your data and the features you switch on. For example, a relatively small file (e.g. 100 variants) **may fail to return results if every variant in the file falls in a different gene and those genes are spread across many chromosomes**. Contrastingly, a file containing 500 variants may return results quickly if those variants all fall in just a few genes.

To mitigate this issue, users should consider **splitting up their input by chromosome** and uploading each chromosome's variants as a separate file. The problem can also be solved by using the VEP script - it is a command line tool, but not as hard to use as you might think! It also offers many more features than the web interface and is a generally much more powerful tool.

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

McLaren W, Pritchard B, Rios D, Chen Y, Flicek P, Cunningham F. *Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor*. **BMC Bioinformatics** 26(16):2069-70(2010) Pubmed PMID: 20562413

Variant Effect Predictor (VEP)

```
1 use strict;
2 use warnings;
3 use Bio::Ensembl::Registry;
4
5 my $registry = 'Bio::Ensembl::Registry';
6
7 $registry->load_registry_from_db(
8     -host => 'ensembl.org',
9     -user => 'anonymous'
10 );
11
12 my $stable_id = 'ENST00000393489'; #this is the stable_id of a human transcript
13 my $transcript_adaptor = $registry->get_adaptor('homo_sapiens', 'core', 'transcript'); #get the adaptor to get the Transcript from the database
14 my $transcript = $transcript_adaptor->fetch_by_stable_id($stable_id); #get the Transcript object
15
16 my $trv_adaptor = $registry->get_adaptor('homo_sapiens', 'variation', 'transcriptvariation'); #get the adaptor to get TranscriptVariation objects
17 my $trvs = $trv_adaptor->fetch_all_by_Transcripts([$transcript]); #get ALL effects of Variations in the Transcript
18
19 foreach my $tv (@{$trvs}) {
20     my $tvas = $tv->get_all_alternate_TranscriptVariationAlleles();
21
22     foreach my $tva(@{$tvas}) {
23         my @ensembl_consequences;
24         my @so_consequences;
25
26         my $ocs = $tva->get_all_OverlapConsequences();
27
28         foreach my $oc(@{$ocs}) {
29             push @ensembl_consequences, $oc->display_term;
30             push @so_consequences, $oc->SO_term;
31         }
32
33         my $sift = $tva->sift_prediction;
34         my $polyphen = $tva->polyphen_prediction;
35
36         print
37             "Variation ", $tv->variation_feature->variation_name,
38             " allele ", $tva->variation_feature_seq,
39             " has consequence ", join(", ", @ensembl_consequences),
40             " (SO ", join(", ", @so_consequences), ").";
41
42         if(defined($sift)) {
43             print " SIFT=$sift";
44         }
45         if(defined($polyphen)) {
46             print " PolyPhen=$polyphen";
47         }
48
49         print "\n";
50     }
51 }
```

EXAMPLE of API usage: Getting all variations in a particular human transcript and see what is the effect of that variation in the transcript

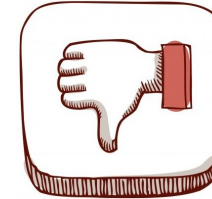
McLaren W, Pritchard B, Rios D, Chen Y, Flicek P, Cunningham F. *Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor*. **BMC Bioinformatics** 26(16):2069-70(2010) Pubmed PMID: 20562413

HPG VARIANT (aka VARIANT)

HPG-VARIANT web site: <http://www.opencb.org/projects/hpg/doku.php?id=variant:overview>



- Free and open source. Part of the [OpenCB Project](#)
- **3 ways of functionality:** C CLI program, Web application and Java RESTful WS API
- **Cloud** variant annotator. Requires **no installation or updates**
- **Regulatory information:** Conserved genomic regions, TFBSs and miRNA targets. ENCODE DNase I hypersensitive sites and Histone methylations
- dbSNP and 1000genomes information
- **Phenotypic information:** HGMD, COSMIC and OMIM
- **Cross-link** with many other DDBB (Ensembl, UniProt, PDB, etc)
- **Input:** VCF, GFF and BED. Accepts compressed files in *tar.gz*
- **11 species** (human, mouse, rat, zebra fish, worm, fly, yeast, dog, pig, mosquito and plasmodium)
- HPG-VARIANT-GWAS to test for **association**



- **Young** program, many **new features coming** in version 2.0 in March 2013
 - PolyPhen and SIFT
 - PhastCons, GERP
 - Many more species (~25 new species)
 - Large structural variants annotation

Medina I, De Maria A, Bleda M, Salavert F, Alonso R, Gonzalez CY, Dopazo J. *VARIANT: Command Line, Web service and Web interface for fast and accurate functional characterization of variants found by Next-Generation Sequencing. Nucleic Acids Research.* 2012 Jul;40(Web Server issue):W54-8 Pubmed PMID: 22693211

HPG VARIANT

Data sources

Core features: genes, transcripts, exons, proteins (UniProt), etc.

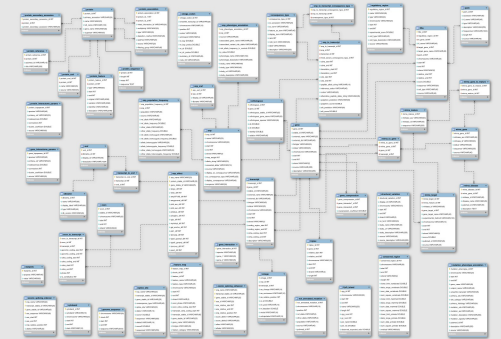
Regulatory: TFBSs, miRNAs, regulatory regions, PWMs, conserved regions, etc.

Functional annotation: OBO ontologies (Gene ontology, disease ontology, etc.), InterPro, etc.

Variation: dbSNP, HapMap, 1000 Genomes project, COSMIC, protein variants, etc.

Systems biology: IntAct, Reactome, gene co-expression, etc.

MySQL cluster



Java RESTful Web Services API design

Structure `ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/{category}/{subcategory}/id/{resource}?{filters}`

Categories **genomic** `ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/genomic/{subcategory}/id/{resource}`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/genomic/region/13:32972105-32973105/snp`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/genomic/variant/13:32972105:A/`

feature `ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/feature/{subcategory}/id/{resource}`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/feature/gene/BRCA2,BCL2/transcript`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/feature/id/BRCA2/xref?dbname=go`

regulatory `ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/regulatory/{subcategory}/id/{resource}`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/regulatory/tf/USF1/tfbs`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/regulatory/mirna_gene/hsa-mir-149/disease`

network `ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/network/{subcategory}/id/{resource}`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/network/pathway/list`
Example: `ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/network/pathway/Triacylglycerol%20biosynthesis/image`

3 ways of functionality:

▶ Command Line Interface (CLI)

[CLI Tutorial](#)

▶ RESTful Web Service

[WS Tutorial](#)

▶ Web Interface

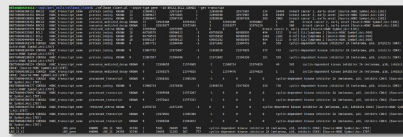
[Web Tutorial](#)

Server
Client

TXT or JSON

Programmatic access


CLI client has been implemented



Web browser access



Usage in web applications



HPG VARIANT - RESTful WS API

Java RESTful Web Services API

General Structure:

`ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/{category}/{subcategory}/id/{resource}`

`ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/genomic/{subcategory}/id/{resource}`

`ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/genomic/variant/id/{resource}`

`ws.bioinfo.cipf.es/cellbase/rest/{version}/{species}/genomic/variant/id/consequence_type`

ID: **Chr:position:reference:variant** (i.e.: 1:150044250:T:G)

`http://ws.bioinfo.cipf.es/cellbase/rest/latest/hsa/genomic/variant/1:150044250:T:G/consequence_type`

HPG VARIANT - RESTful WS API

Column position	Column header	Description
1	chromosome	Chromosome number
2	position	Position of the variant
3	reference	Reference allele at this position
4	alternative	Non-reference allele called at this position
5	feature ID	Ensembl identifier or name of a genomic feature where the variant has been found. In the case of introns and intergenic regions this field is empty
6	feature name	HGNC symbol of the feature (if exists)
7	consequence type	Consequence type code (see more here)
8	biotype	Biotype of the target transcript
9	feature chromosome	Chromosome number for the feature
10	feature start	Feature start position
11	feature end	Feature end position
12	feature strand	Feature strand
13	snp ID	Single Nucleotide Polymorphism (SNP) rsID. This field is empty unless the position is annotated as a SNP
14	ancestral allele	Ancestral allele
15	alternative allele	Variant allele found
16	gene Ensembl ID	Gene Ensembl ID related with this feature
17	Ensembl transcript ID	Transcript Ensembl ID related with this feature
18	gene name	Gene name related with this feature
19	SO consequence type ID	Sequence Ontology consequence type ID
20	SO consequence type name	Sequence Ontology consequence type name
21	consequence type description	Consequence type description
22	consequence type category	Consequence type category as specified here
23	aminoacid position	If the variant is in a coding region for a transcript it indicates the affected aminoacid position
24	aminoacid change	If the variant is non-synonymous for a transcript it indicates the original and the new aminoacids in 1-letter code
25	codon change	If the variant is non-synonymous for a transcript it indicates the original and the new codons

HPG VARIANT - CLI program

- ▶ Create a folder called `effect` in your `mda13` directory:

```
mkdir /home/biouser/mda13/effect
```

- ▶ Download the program and save it into your `mda13/effect` folder :

<http://www.opencb.org/projects/hpg/doku.php?id=variant:downloads>

- ▶ Extract the contents

- ▶ Download the example file, save it into the `mda13/effect/hpg-variant-bin-0.4/` folder and extract the content.







http://bioinfo.cipf.es/courses/mda13genomics/program?&#variant_annotation

- ▶ Open a Terminal  and move to the folder where files have been extracted:

```
cd /home/biouser/mda13/effect/hpg-variant-bin-0.4/
```

- ▶ Execute it!

```
hpg-var-effect -v CHB.exon.2010_03.sites.vcf --outdir effect_output/
```

	Binaries	Sources
Debian 6	 Binary package	
Ubuntu 12.04	 Binary package	
Fedora 17	 Binary package	 Source package
Other	 Zipped binaries *	 Source tarball

* Only for Debian 6 / Ubuntu 10.04 or greater

HPG VARIANT - CLI program

▶ Open the `effect_output` directory and let's see what is in!

▶ **Is there any variant affecting a splice donor site?**

Yes! There are 4 variants in a splice donor site affecting 5 genes: CEP104, HSE4, C1orf159, CALML6 and C1ORF222.

▶ **Is there any SNP with known phenotype?**

Yes! There are 39 SNPs with known phenotype. They have been associated with Migraine, Alzheimer's disease, Ulcerative colitis, Rheumatoid arthritis, etc.

HPG VARIANT - Web application

► Web application

<http://variant.bioinfo.cipf.es/>

The screenshot shows the web application interface for the Variant analysis tool. At the top right, there is a "sign in" button. Below the header, the page title is "Variant analysis tool" with sub-links for "Variant effect" and "VCF Viewer". A navigation bar includes "Home", "home", "help", "tutorial", "Projects", and "Upload data".

Overview

VARIANT (VARIant ANalysis Tool) can report the functional properties of any variant in all the human, mouse or rat genes (and soon new model organisms will be added) and the corresponding neighborhoods. Also other non-coding extra-genic regions, such as miRNAs are included in the analysis.

VARIANT not only reports the obvious functional effects in the coding regions but also analyzes noncoding SNVs situated both within the gene and in the neighborhood that could affect different regulatory motifs, splicing signals, and other structural elements. These include: Jaspas regulatory motifs, miRNA targets, splice sites, exonic splicing silencers, calculations of selective pressures on the particular polymorphic positions, etc.

Note

This web application makes an intensive use of new web technologies and standards like HTML5, so browsers that are fully supported for this site are: Chrome 14+, Firefox 7+, Safari 5+ and Opera 11+. Older browser like Chrome13-, Firefox 5- or Internet Explorer 9 may rise some errors. Internet Explorer 6 and 7 are no supported at all.

Sign in

You must be logged in to use this Web application, you can **register** or use a **anonymous user** as shown in the following image by clicking on the **"Sign in"** button on the top bar

The image shows a "Sign in" dialog box with the following elements:

- Input fields for "e-mail:" and "password:"
- A checked checkbox for "Anonymous login" with the note "Your work will be lost after logout session".
- Radio buttons for "Anonymous selected" and "No password required".
- Buttons for "Sign in", "Forgot your password?", and "New account".

The "sign in" button in the top navigation bar and the "Sign in" button in the dialog box are highlighted with red boxes.

Thank you!



Questions?