# Database Integration

Paul Flicek

Vertebrate Genomics

EMBL-EBI

# (Dramatically) Simplified Clinical Workflow

Identify variants

Technically easy and getting easier

Use what we already know to make some sense of them

...nething ...but it

EMBL-EBI

# Data interpretation: beyond research toward medical practice

- Needs:
  - Consistent, traceable data generation and analysis routines
  - Robust annotation based on public information sources such as those at the EBI and NCBI
    - Probably 95% of all information that could be used to understand and interpret human variation is already in the public domain
  - Reporting into medical records

EMBL-EBI

# Database integration

- Part 1: Continually update the existing information to ensure it is accurate and comprehensive

- Part 2: Provide some method to search relevant resources using variants and/or whole genomes as input

EMBL-EBI

*european genome-phenome archive*

- EGA Home
- :: Submit to EGA
- :: Information
- :: Help
- Contact Us

**EXPLORE THE EGA**

- Browse studies
- Browse datasets
- Browse data access committees
- Browse data providers

**USER LOGIN**

Username: *

Password: *

Request new password    Log in

## The European Genome-phenome Archive (EGA)

The European Genome-phenome Archive (EGA) repository allows you to **explore datasets** from numerous genotype experiments, supplied by a range of **data providers**.

### Studies

Studies are experimental investigations of a particular phenomenon or trait.

**Browse all studies**

### Datasets

The EGA archives a large number of datasets, the access to which is controlled by a dataset access committee (DAC).

**Browse datasets that we hold**

**Browse the list of DACs**

### Data Providers

Data Providers can be involved in creating studies, data submission and the designation of data access committees (DACs).

**Browse EGA Data Providers**

### Learn about the EGA

- Introduction to the EGA
- How to obtain an account with the EGA

**Video resources**

Watch a video of how to use your EGA account

### Help

- Users FAQ
- Submitters FAQ
- FTP & Aspera FAQ
- Contact Us
- EGA mailing list

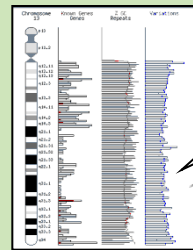# The European Genome-phenome Archive

- Secure storage and authorised access to all types of data sets that might be generated in the context of research into molecular medicine
  - DNA sequence; Array-based genotypes; epigenetic data
  - Transcriptomics; Proteomics
  - Phenotype data

- Used for GWAS, ICGC, IHEC, IHMC, UK10K and data

- EGA supports only data access decisions that are based on original consent
  - Authorized users have personal accounts in our system
  - Access to the data requires account password
  - Data decryption requires a separate key that must be requested and is sent off line

EMBL-EBI
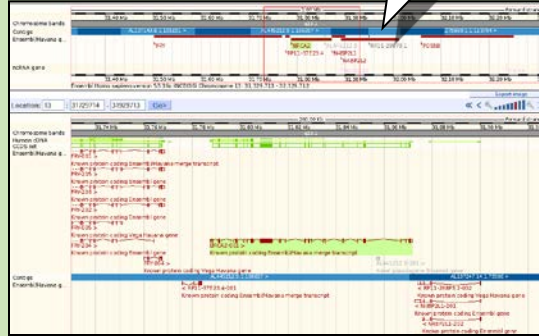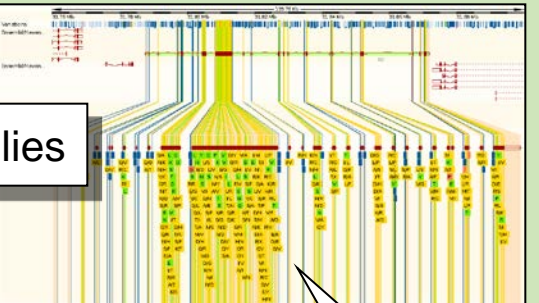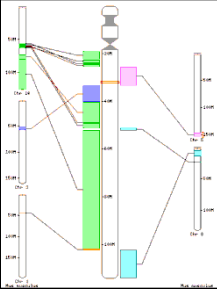
# Ensembl genome-wide annotation

Genomic alignments
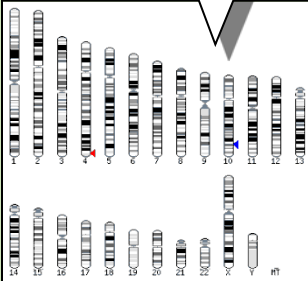
Chromosomes
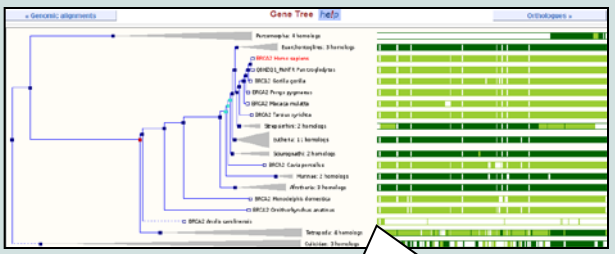
Genes

Pick a genome

Synteny

Gene families

SNPs

Across species

Orthology

Within species

# Integrating variation data across the genome

- Polymorphism data (from dbSNP)
  - SNPs and indels for 14 species including 1000 Genomes
  - Allele and genotype frequencies by population
  - Locus-specific data from LRG
- Structural polymorphism data
- Mutation data (human)
  - Somatic mutation data (from COSMIC)
  - Human Gene Mutation Database (HGMD) IDs

- Phenotype associations: OMIM, UniProt, GWAS
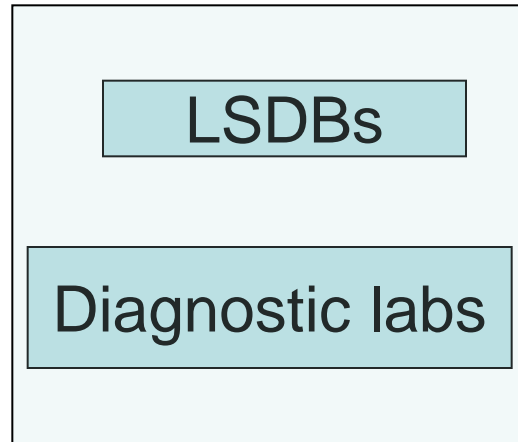- Affymetrix and Illumina chipsets

# Variation annotation – phenotype data

- 37,964 somatic mutations:
  - COSMIC
- 57,930 germline mutations:
  - HGMD
- 56,177 literature curation:
  - OMIM
  - UniProt
- 62,737 GWAS data:
  - NHGRI GWAS catalog
  - Open Access DB
  - EGA
- 22,449 from SNPedia by DAS

# Variation annotation – phenotype data

LSDBs

Diagnostic labs

Locus-specific information



Dalgleish, et al.
*Genome Medicine*
2010



Genome-wide information

- LRG project- Locus Reference Genomic
  - Create stable reference sequences (LRGs)
  - Use LRGs for exchange of variation data

wellcome trust
sanger
institute

e!

EMBL-EBI

# Database integration

- Part 1: Continually update the existing information to ensure it is accurate and comprehensive

- Part 2: Provide some method to search relevant resources using variants and/or whole genomes as input

EMBL-EBI

# Ensembl Variant Effect Prediction (VEP) tool

- Calculates the effect of SNPs in the context of Ensembl genes and regulatory features
  - Web and API interface
  - Code back-ported to support NCBI36 assembly
  - Programmatic support for tab-delimited and VCF files
  - Easily integrated into analysis pipelines
- Working within ICGC to capture structural and other genome rearrangements
- Disruption of experimentally observed TF binding sites and conserved regions
- Ability to run without connection to the internet
- Support for user defined analysis plug-ins coming in January 2012
- Will return if variant is present in EGA dataset in 2012
- Effectively a variant based search of EBI's data resources

McLaren, et al. *Bioinformatics.* 2010

EMBL-EBI

# Ensembl VEP Implementation

**Input file**

Species: Human (Homo sapiens): GRCh37

50+ species at www.ensembl.org
300+ at www.ensemblgenomes.org

Name for this upload (optional):

Paste file:

Data input by file upload or external URL

Upload file: Choose File No file chosen

or provide file URL:

Support for multiple file formats: VCF, Pileup, HGVS, dbSNP rsID

Input file format: Ensembl default

**Options**

Get regulatory region consequences (human and mouse only): ☑

Type of consequences to display: Ensembl terms

Output Ensembl, Sequence Ontology (SO) or NCBI consequence terms

Check for existing co-located variants: Yes

Return results for variants in coding regions only: ☐

Show HGNC identifier for genes where available: ☐

Find existing overlapping variants annotated by Ensembl

Show Ensembl protein identifiers where available: ☐

Show HGVS identifiers for variants where available: No

**Non-synonymous SNP predictions (human only)**

Create HGVS notations

SIFT predictions: No

PolyPhen predictions: No

Condel consensus (SIFT/PolyPhen) predictions: No

**Frequency filtering of existing variants (human only)**

Filter variants by frequency: ☐

Include **SIFT**, **PolyPhen** and **Condel** predictions for non-synonymous changes in human

NB: Enabling frequency filtering may be very slow for large datasets

Filter: Exclude variants with MAF greater than 0.1 in any 1KG low coverage population

Filter input against **HapMap** or **1000 genomes** frequency data

Next >

# Output

| Uploaded Variation | Location | ...ence | Position in cDNA | Position in CDS | Position in protein | Amino acid change | Codon change | Co-located Variation | Extra |
|---|---|---|---|---|---|---|---|---|---|
| rs10576 | 21:2696517... | ...s_codon | 915 | 873 | 291 | P | ccA/ccG | rs10576 | **HGNC**=MRPL39 |
| rs10576 | 21:2696517... | ...s_codon | 852 | 843 | 281 | P | ccA/ccG | rs10576 | **HGNC**=MRPL39 |
| rs10576 | 21:2696517... | ...s_codon | 915 | 873 | 291 | P | ccA/ccG | rs10576 | **HGNC**=MRPL39 |
| rs1057885 | 21:2696520... | ...s_codon | 882 | 840 | 280 | V | gtA/gtG | rs1057885 | **HGNC**=MRPL39 |
| rs1057885 | 21:2696520... | ...s_codon | 882 | 840 | 280 | V | gtA/gtG | rs1057885 | **HGNC**=MRPL39 |
| rs1057885 | 21:2696520... | ...s_codon | 819 | 810 | 270 | V | gtA/gtG | rs1057885 | **HGNC**=MRPL39 |
| rs113417859 | 21:4019154... | ...ion_variant | - | - | - | - | - | rs113417859 | - |
| rs113417859 | 21:4019154... | ...pt_variant, ...ariant | - | - | - | - | - | rs113417859 | **HGNC**=TMPRSS3 |
| rs113417859 | 21:4019154... | ...s_codon | 1393 | 933 | 311 | F | ttC/ttT | rs113417859 | **HGNC**=ETS2 |
| rs113417859 | 21:4019154... | ...s_codon | 1223 | 933 | 311 | F | ttC/ttT | rs113417859 | **HGNC**=ETS2 |
| rs113417859 | 21:4019154... | ...eam_variant | - | - | - | - | - | rs113417859 | **HGNC**=ETS2 |
| rs113417859 | 21:4019154... | ...am_variant | - | - | - | - | - | rs113417859 | **HGNC**=ETS2 |
| rs1135638 | 21:2696514... | ...s_codon | 939 | 897 | 299 | G | ggC/ggT | rs1135638 | **HGNC**=MRPL39 |
| rs1135638 | 21:2696514... | ...s_codon | 876 | 867 | 289 | G | ggC/ggT | rs1135638 | **HGNC**=MRPL39 |
| rs1135638 | 21:2696514... | ...s_codon | 939 | 897 | 299 | G | ggC/ggT | rs1135638 | **HGNC**=MRPL39 |
| rs114053718 | 21:3402919... | ...ous_codon | 2722 | 2597 | 866 | I/T | aTt/aCt | rs114053718 | **SIFT**=deleterious(0); **HGNC**=SYNJ1 |
| rs114053718 | 21:3402919... | ...ous_codon | 2714 | 2714 | 905 | I/T | aTt/aCt | rs114053718 | **SIFT**=deleterious(0); **HGNC**=SYNJ1 |
| rs114053718 | 21:3402919... | ...ous_codon | 2597 | 2597 | 866 | I/T | aTt/aCt | rs114053718 | **SIFT**=deleterious(0); **HGNC**=SYNJ1 |
| rs114053718 | 21:3402919... | ...t_variant | 384 | - | - | - | - | rs114053718 | **HGNC**=SYNJ1 |
| rs114053718 | 21:3402919... | ...ous_codon | 2722 | 2714 | 905 | I/T | aTt/aCt | rs114053718 | **SIFT**=deleterious(0); **HGNC**=SYNJ1 |

Show All entries | ...hide columns | Filter

# Sequence Ontology consequences

- Provides a structured controlled vocabulary for the description of mutations at both the sequence and more gross level in the context of genomic databases



- regulatory region
- TF binding site

- intergenic
- 5KB upstream
- 2KB upstream

- 5 prime UTR
- initiator codon change

- synonymous codon
- non-synonymous codon
- inframe codon gain
- inframe codon loss
- stop gained
- frameshift
- coding sequence variant
- complex change in transcript

- splice donor
- splice acceptor

- splice region
- intron

- stop lost
- stop retained
- incomplete terminal codon

- 3 prime UTR

- 500B downstream
- 2KB downstream
- 5KB downstream

# SIFT, PolyPhen and Condel in practice

- Store every possible score for every* protein

| | A | C | D | E | F | G | H | I | K | L | M | N | P | Q | R | S | T | V | W | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 0.001 | 0.047 | 0.007 | 0.007 | 0.007 | 0.002 | 0.047 | 0.001 | 0.002 | 0.001 | - | 0.007 | 0.007 | 0.007 | 0.007 | 0.002 | 0.002 | 0.001 | 0.094 | 0.017 |
| **2** | 0.081 | 0.547 | 0.547 | 0.348 | 0.201 | 0.348 | 0.817 | 0.081 | 0.348 | - | 0.348 | 0.547 | 0.547 | 0.547 | 0.547 | 0.201 | 0.201 | 0.081 | 0.817 | 0.547 |
| **3** | 0.007 | 0.191 | 0.007 | 0.002 | 0.094 | 0.017 | 0.094 | 0.047 | 0.002 | 0.017 | 0.094 | 0.017 | 0.017 | - | 0.007 | 0.007 | 0.017 | 0.017 | 0.191 | 0.047 |
| **4** | 0.017 | 0.362 | 0.201 | 0.106 | 0.106 | 0.106 | 0.362 | 0.017 | 0.106 | 0.017 | 0.201 | 0.362 | 0.201 | 0.362 | 0.362 | 0.106 | 0.04 | - | 0.677 | 0.201 |
| **5** | 0.017 | 0.362 | 0.201 | 0.106 | 0.106 | 0.106 | 0.362 | 0.017 | 0.106 | 0.017 | 0.201 | 0.362 | 0.201 | 0.362 | 0.362 | 0.106 | 0.04 | - | 0.677 | 0.201 |
| **6** | 0.007 | 0.191 | 0.007 | 0.002 | 0.094 | 0.017 | 0.094 | 0.047 | 0.002 | 0.017 | 0.094 | 0.017 | 0.017 | - | 0.007 | 0.007 | 0.017 | 0.017 | 0.191 | 0.047 |
| **7** | 0.081 | 0.817 | 0.035 | - | 0.547 | 0.081 | 0.547 | 0.547 | 0.081 | 0.201 | 0.547 | 0.201 | 0.201 | 0.081 | 0.201 | 0.081 | 0.081 | 0.201 | 0.817 | 0.547 |
| **8** | 0.663 | 0.99 | 0.964 | 0.964 | 0.964 | - | 0.99 | 0.964 | 0.964 | 0.964 | 0.99 | 0.922 | 0.964 | 0.964 | 0.964 | 0.848 | 0.964 | 0.964 | 0.99 | 0.99 |
| **9** | 0.081 | 0.817 | 0.081 | 0.081 | 0.547 | 0.081 | 0.348 | 0.547 | 0.081 | 0.201 | 0.547 | - | 0.348 | 0.201 | 0.201 | 0.081 | 0.081 | 0.201 | 0.817 | 0.547 |
| **...** | | | | | | | | | | | | | | | | | | | | |

- Condel scores are an algorithmic function of SIFT and Polyphen scores

# Regulatory region consequences

- Variant within a regulatory feature = RegulatoryFeature

- Variant within a transcription factor binding motif = MotifFeature

- Variant in an "informative position" = HIGH_INF_POS



| rs75265131 | 12:96156035 | C | - | | MA0074.1 | MotifFeature | REGULATORY_REGION | rs75265131 | **MATRIX=**Jaspar_Matrix_RXRA::VDR:MA0074.1; **HIGH_INF_POS=**Y |
|---|---|---|---|---|---|---|---|---|---|
| rs75265131 | 12:96156035 | C | ENSG00000074527 | ENST00000547980 | | Transcript | INTRONIC | rs75265131 | - |
| rs75265131 | 12:96156035 | C | - | | ENSR00000435320 | RegulatoryFeature | REGULATORY_REGION | rs75265131 | - |

# Has this variant ever been seen before?

- Quickly becoming the most common question in human genomics
  - Incredibly hard to answer
- Nature said (in the October 2010 1000 Genomes issue) that about 2700 genomes had been sequenced and estimate 30,000 by the end of 2011
  - Beyond the those currently in the 1000 Genomes project (~2000)relatively few of these genomes are easily accessible
- There are many more exomes
  - Access here can be a problem as well

- Some data is available under controlled access and the fraction of data in this category is expected to increase

# Future

- Ensembl is not a clinical decision support tool and only a fraction of the important resources were presented

- It does show the way forward
  - Comprehensive
  - Versioned
  - Standardized
  - Using controlled terminology
  - Regularly updated
  - Evidence based and algorithmic
  - Fully open

- There is uncertainty at every step in the process from the genome reference to the gene set to the interpretation and we have to work in this environment

EMBL-EBI

# Acknowledgements

- Ensembl Annotation and VEP: Will McLaren, Graham Ritchie, Pontus Larsson, Daniel Sobral, Bethan Yates, Anne Parker, Jackie MacArthur, Fiona Cunningham

- EBI Variation Archives: Ilkka Lappalainen, Vasudev Kumanduri, Dylan Spalding, Mick Maguire, Lisa Skipper, Jeff Almeida-King

- Funding: Wellcome Trust, European Commission, NHGRI, British Heart Foundation, EMBL

EMBL-EBI

# EBI data integration and added value

- EBI search provides integration into EBI existing spines (DAS based)
- Development of new spines diseases, cell type, tissue, tools
- User focussed design with general and specific user groups
- Added value  - terminology, literature searching, pathways etc (user defined)
- Reciprocal integration between KOMP2 web portal and EBI resources



EMBL-EBI

### Alkaline phosphatase, liver/bone/kidney Gene

## Alpl differential expression summary

☑ View in Gene B

### Organism part

liver, kidney, placenta, thymus
☑ View all

### Disease state

normal, (empty), Pb-A Infected, myocardial infarction
☑ View all

### Cell type

embryonic stem cell, hematopoietic stem cell, T effector cell, T regulatory cell
☑ View all

### Cell line

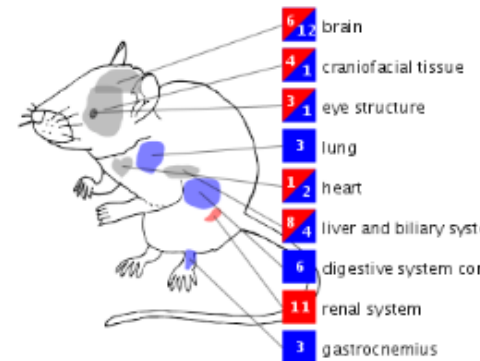N1E-115 wild_type, 67NR, Swiss5, Swiss8
☑ View all

### Compound treatment

control, none, vehicle - arachis oil, 17beta-estradiol
☑ View all

### Developmental stage

embryo, adult, fetus, neonate
☑ View all

Gene

Expression

Protein

Protein Structure

Literature

brain
craniofacial tissue
eye structure
lung
heart
liver and biliary syst
digestive system co
renal system
gastrocnemius

☑ Number of published studies where the gene over/unde
compared to the gene's overall mean expression level i

**House Mouse**
*Mus musculus*

Gene → KOMP2 Ensembl links

Expression → LacZ summaries, image links

Protein

Protein Structure

Literature

Disease → Mouse models of disease, phenotype summaries

Pathways

Tissues → Expression summaries, phenotype links

Chemistry

Tools → Mouse knockouts, phenotype summaries, CDA links