# Module 7

# Variation, Function and Disease

**Bert Overduin, Ph.D.**

Edinburgh Genomics
The University of Edinburgh
Edinburgh EH9 3FL
Scotland, United Kingdom

# Overview

- Tools to predict variant effects
- Gene expression databases
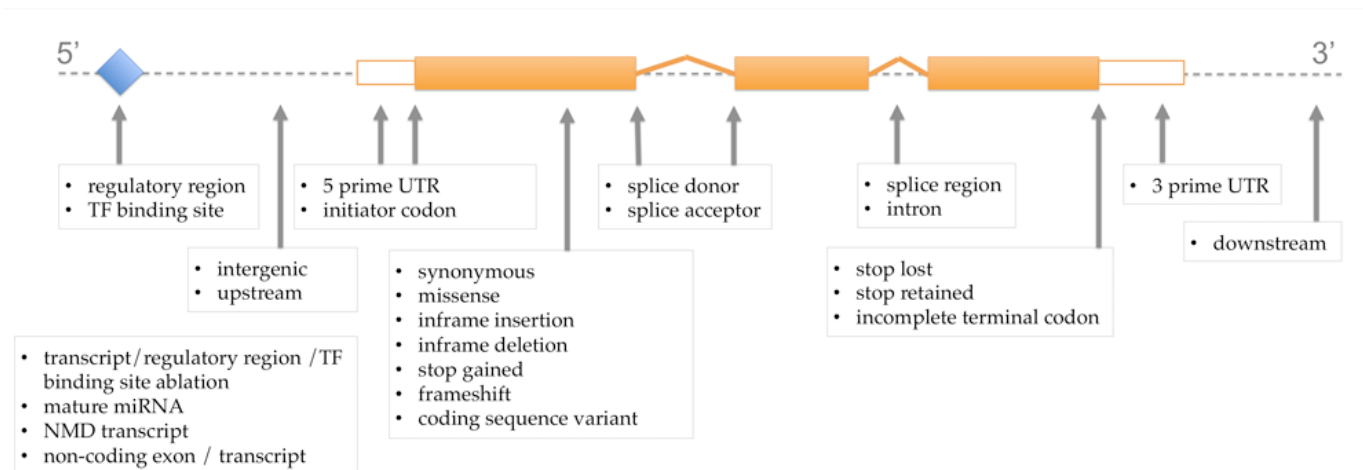- "Disease" / phenotype databases
- Ontologies

# Variant Effects

Tools to categorise and prioritise newly discovered variants:

- Variant Effect Predictor (Ensembl)
- Variation Annotation Integrator (UCSC)
- PolyPhen-2
- SIFT

# Variant Effect Predictor (VEP)

- Predicts functional consequences of known and unknown variants

- For substitutions, insertions, deletions and structural variants

# VEP Output

- Affected genes / transcripts / regulatory features / motifs
- Gene symbols
- IDs from Ensembl, CCDS, UniProt, HGVS
- Consequence (e.g. missense, stop gained, stop lost)
- Location of variant
- Co-located known variant(s)
- Minor allele frequencies from the 1000 Genomes Project
- PolyPhen and SIFT prediction and score

etc. etc.

# Polyphen-2 and SIFT

- Predict the effect of missense variants

**PolyPhen-2** (<u>Poly</u>morphism <u>Phen</u>otyping)

- Uses physical and comparative considerations
- Scale from 0 (benign), via possibly damaging to 1 (probably damaging)

**SIFT** (<u>S</u>orting <u>I</u>ntolerant <u>F</u>rom <u>T</u>olerant)

- Uses the degree of conservation of amino acid residues
- Scale from 0 (deleterious) to 1 (tolerated)

# Warning

- All these tools do is make predictions!
- Findings should always be confirmed by doing experiments!

# Gene Expression Databases

**GEO Profiles** (NCBI)

- Gene expression profiles

- Derived from GEO (<u>G</u>ene <u>E</u>xpression <u>O</u>mnibus)

**Expression Atlas** (EBI)

- Baseline Atlas: shows which gene products are present (and at what abundance) in "normal" conditions

- Differential Atlas: shows genes that are up- or down-regulated in a wide variety of different experimental conditions

- Derived from ArrayExpress

# Baseline Atlas

# Differential Atlas

# Disease / Phenotype Databases

**OMIM**

- <u>O</u>nline <u>M</u>endelian <u>I</u>nheritance in <u>M</u>an
- Catalog of human genes and genetic disorders

**COSMIC**

- <u>C</u>atalog <u>O</u>f <u>S</u>omatic <u>M</u>utations <u>I</u>n <u>C</u>ancer
- Database of somatically acquired mutations found in human cancer

**DECIPHER**

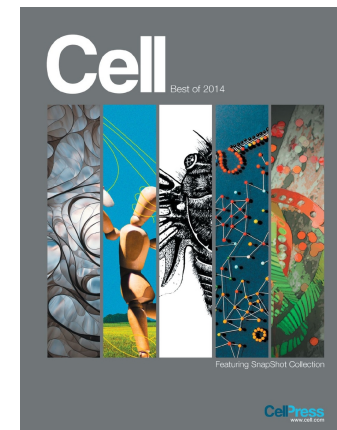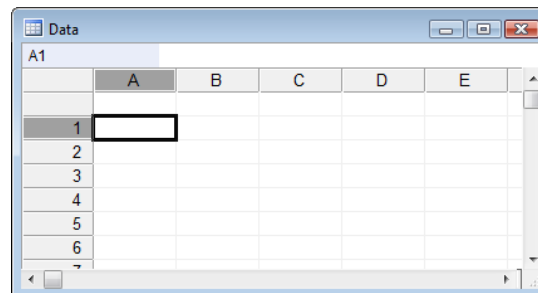- <u>D</u>atabas<u>E</u> of genomi<u>C</u> var<u>I</u>ation and <u>P</u>henotype in <u>H</u>umans using <u>E</u>nsembl <u>R</u>esources)
- Database of genomic variation data from analysis of patient DNA

# Ontologies

What do you think about when you hear the word "cell"?

# Ontologies

What do you think about when you hear the word "cell"?

# Ontologies

**Gene Ontology** (GO)

- Describes gene products in terms of their:
  - Associated biological processes (what?)
  - Cellular components (where?)
  - Molecular functions (how?)

**Sequence Ontology** (SO)

- Describes features and attributes of biological sequence

# Worked examples

# Worked Examples URLs

- Ensembl Variant Effect Predictor: http://www.ensembl.org/info/docs/tools/vep/index.html
- UCSC Variation Annotation Integrator: http://genome.ucsc.edu/cgi-bin/hgVai
- PolyPhen-2: http://genetics.bwh.harvard.edu/pph2
- SIFT: http://sift.jcvi.org
- OMIM: http://omim.org
- GEO Profiles: http://www.ncbi.nlm.nih.gov/geoprofiles
- Expression Atlas: http://www.ebi.ac.uk/gxa/home
- COSMIC: http://cancer.sanger.ac.uk/cosmic/
- DECIPHER: https://decipher.sanger.ac.uk/
- Gene Ontology: http://geneontology.org/
- Sequence Ontology: http://www.sequenceontology.org/