# Module 7: Variation, Function and Disease

**Aim**

Learn how to explore variation and the relationship between genotype and phenotype / disease using the following tools and databases:

- The Ensembl Variation Effect Predictor (VEP)
- PolyPhen-2
- OMIM
- GEO Profiles and Gene Expression Atlas
- COSMIC
- DECIPHER
- Ontologies

Often the most valuable information to know about a variant is the effect the observed alleles have on genes, transcripts and proteins. This information can be very helpful to prioritise any variants for further investigation. To determine this effect, several tools are available. One should keep in mind though that all these tools do is make predictions and consequently findings should always be confirmed by experiments.

**The Ensembl Variant Effect Predictor (VEP)**

The Ensembl Variant Effect Predictor (VEP) determines the effect of variants (SNPs, insertions, deletions, CNVs or structural variants) on genes, transcripts, and protein sequence, as well as regulatory regions. It also calculates SIFT and PolyPhen scores for changes to protein sequence.

**Worked example 1: The Ensembl Variant Effect Predictor (VEP)**

In this worked example we will study four newly found variants in human:

Deletion of an A at position 128328461 on chromosome 9
Substitution C>A at position 128322349 on chromosome 9
Substitution C>G at position 128323079 on chromosome 9
Substitution G>A at position 128322917 on chromosome 9

We will use the **Ensembl VEP** to answer the following questions:

- Have my variants already been annotated in Ensembl?

- What genes are affected by my variants?
- Do my variants result in protein changes?
- Do any of my variants affect gene regulation?

Go to the Ensembl Variant Effect Predictor page (http://www.ensembl.org/info/docs/tools/vep/index.html).

This page contains information about the VEP, including links to download the script version of the tool. Click on "Launch Ve!P" to open the input form.



The data should be inputted in the following format:

Chromosome Start End Alleles (reference/mutation) Strand Name

Replace the example data in the "Paste data" box with:

```
9 128328461 128328461 A/- + var1
9 128322349 128322349 C/A + var2
9 128323079 128323079 C/G + var3
9 128322917 128322917 G/A + var4
```

The VEP will automatically detect that the data is in "Ensembl default format".

There are further options that you can choose for your output. These are categorised as "Identifiers and frequency data", "Filtering options" and "Extra options". Let's open all the menus and take a look.

Hovering with your mouse over an option will show a pop-up with an explanation of that option.

When you've selected everything you need, scroll right to the bottom and click [Run].



**Your ticket number**

**Click to get your results**

**Buttons to save, edit or delete your job**

The display will show you the status of your job. It will say "Queued", then automatically switch to "Done" when the job is done, you do not need to refresh the page. You can edit or discard your job at this time. If you have submitted multiple jobs, they will all appear here.

Click "[View results]" once your job is done.

In your results you will see a graphical summary of your data, as well as a table of your results. (Note that some empty columns in the results table have been hidden in the following screenshot to save space.)

Note that the UCSC Genome Browser has a similar tool, named the Variation Annotation Integrator (http://genome.ucsc.edu/cgi-bin/hgVai), that offers largely the same functionality as the Ensembl Variant Effect Predictor.

**PolyPhen-2**

PolyPhen-2 (Polymorphism Phenotyping v2) is a tool that predicts the possible impact of an amino acid substitution on the structure and function of a human protein using physical and comparative considerations.

**Worked example 2: PolyPhen-2**

In this worked example we will have a look at a variant of the human *BRAF* (B-Raf proto-oncogene, serine/threonine kinase) gene.  The most common variant in the BRAF protein is a V to E change at amino acid position 600. Let's assume we did not know the consequence of this variant and use PolyPhen-2 to see if it is deleterious. The UniProt identifier for the BRAF protein is P15056 (http://www.uniprot.org/uniprot/P15056).

> **STEP 1** – **Go to** the PolyPhen homepage:
> http://genetics.bwh.harvard.edu/pph2



**PolyPhen-2** prediction of functional effects of human nsSNPs

| Home | About | Help | Downloads | Batch query | WHESS.db |

PolyPhen-2 (**Poly**morphism **Phen**otyping v2) is a tool which predicts possible impact of an amino acid substitution on the structure and functio... and comparative considerations. Please, use the form below to submit your query.

**15-Feb-2012:** PolyPhen-2 server has been updated to utilize **version 2.2.2** of the software, protein sequences from UniProtKB/UniRef10... PDB/DSSP Snapshot 03-Jan-2012 (78,304 entries) and **UCSC** MultiZ multiple alignments of 45 vertebrate genomes with hg19/GRCh37 human genome...

**Query Data**

Protein or SNP identifier

Protein sequence in FASTA format

Position

Substitution   AA₁  A R N D C E Q G H I L K M F P S T W Y V
               AA₂  A R N D C E Q G H I L K M F P S T W Y V

Query description

[Submit Query] [Clear] [Check Status]

Display advanced query options

> **STEP 2** – **Enter** P15056 into the 'Protein identifier' textfield

> **STEP 3** – **Put** 600 in the 'Position' textfield

> **STEP 4** – **Select** V for 'AA₁' and E for 'AA₂'

> **STEP 5** – **Click** [Submit Query]

A similar tool to PolyPhen is SIFT (Sorting Intolerant From Tolerant; http://sift.jcvi.org). SIFT prediction is based on the degree of conservation of

amino acid residues in sequence alignments derived from closely related sequences, collected through PSI-BLAST. SIFT can be applied to naturally occurring nonsynonymous polymorphisms or laboratory-induced missense mutations.

In the following section we will have a look at several resources that focus on the relationship between genotype and phenotype / disease.

## OMIM

The OMIM (Online Mendelian Inheritance in Man; http://omim.org) database is a catalogue of human genes and genetic disorders.  OMIM focuses on the relationship between phenotype and genotype.

### Worked example 3: OMIM

In this worked example we will have a look what information OMIM contains about the human *BRAF* gene.

> **STEP 1** – **Go to** the OMIM homepage:
> http://omim.org

BRAF [Search] Sort by: ● Relevance ○ Date updated

Advanced Search: OMIM, Clinical Synopses, OMIM Gene Map   Toggle: search terms highlighted
Search History: View, Clear

Retrieve corresponding: [gene map] [clinical synopses]

Search: 'BRAF'
Results: 1 - 10 of 64 | Show all | 1 2 3 4 5 6 7 Next Last

1 : * 164757. V-RAF MURINE SARCOMA VIRAL ONCOGENE HOMOLOG B1; BRAF          Gene Tests, Links
      BRAF/AKAP9 FUSION GENE, INCLUDED
      Cytogenetic location: 7q34 , Genomic coordinates (GRCh37): 7:140,433,811 - 140,624,
      Matching terms: braf

**STEP 3 – Click on '*164757' to retrieve the gene information**

2 : % 155600. MELANOMA, CUTANEOUS MALIGNANT, SUSCEPT          ICD+, Links
      Cytogenetic location: 1p36 , Genomic coordinates (GRCh37): 1:0 - 28,000,000
      Matching terms: braf

3 : # 115150. CARDIOFACIOCUTANEOUS SYNDROME          Gene Tests, ICD+, Links
      Cytogenetic locations: 7q34 , 12p12.1 , 15q22.31 , 19p13.3
      Matching terms: braf

Gene centric information is preceded by '*'. Disease centric information is preceded by '#' and '%' (if the underlying molecular basis is not known).

4 : # 114500. COLORECTAL CANCER; CRC
      Cytogenetic locations: 1p36.13 , 1p22.3 , 1p13.2 , 2p25.1 , 3q26.32 , 4p16.3 , 4q31.3 , 5q2
      22q13.2
      Matching terms: braf

5 : * 613344. KIAA1549 GENE; KIAA1549
      KIAA1549/BRAF FUSION GENE, INCLUDED
      Cytogenetic location: 7q34 , Genomic coordinates (GRCh37): 7:138,516,125 - 138,666,
      Matching terms: braf

6 : * 604001. A-KINASE ANCHOR PROTEIN 9; AKAP9          Gene Tests, Links
      AKAP9/BRAF FUSION GENE, INCLUDED
      Cytogenetic location: 7q21.2 , Genomic coordinates (GRCh37): 7:91,570,188 - 91,739,986
      Matching terms: braf

7 : # 188550. THYROID CARCINOMA, PAPILLARY          Gene Tests, ICD+, Links
      Cytogenetic locations: 1p13.2 , 7q33-q34 , 8p22 , 10q11.23 , 10q21.2 , 14q32.12 , 17q24.2
      Matching terms: braf

The returned entry is a detailed description of the gene, its location and genetic defects that result in disease. Each piece of clinical and genetic data is cited providing an excellent platform for understanding the function of the gene in a disease setting.

BRAF

Search

Advanced Search ▾ | Search History | Display Options ▾

**\*164757**

**V-RAF MURINE SARCOMA VIRAL ONCOGENE HOMOLOG B1; BRAF**

*Alternative titles; symbols*
ONCOGENE BRAF
BRAF1
RAFB1

Other entities represented in this entry:

**BRAF/AKAP9 FUSION GENE, INCLUDED**

BRAF/KIAA1549 FUSION GENE, INCLUDED

*HGNC Approved Gene Symbol: BRAF*

*Cytogenetic location: 7q34*      *Genomic coordinates (GRCh37): 7:140,415,748-140,624,563* (from NCBI)

**Gene-Phenotype Relationships**

| Location | Phenotype | Phenotype MIM number | Phenotype mapping key |
|----------|-----------|----------------------|------------------------|
| 7q34 | Adenocarcinoma of lung, somatic | 211980 | 3 |
| | Cardiofaciocutaneous syndrome | 115150 | 3 |
| | Colorectal cancer, somatic | | 3 |
| | LEOPARD syndrome 3 | 613707 | 3 |
| | Melanoma, malignant, somatic | | 3 |
| | Nonsmall cell lung cancer, somatic | | 3 |
| | Noonan syndrome 7 | 613706 | 3 |

**Table of Contents for \*164757**
Title
Gene-Phenotype Relationships
Text
  Cloning and Expression
  Gene Function
  Mapping
  Molecular Genetics
  Cytogenetics
  Animal Model
Allelic Variants
  Table View
References
Contributors
Creation Date
Edit History
**External Links for Entry:**
▸ Genome
▸ DNA
▸ Protein
▸ Gene Info
▸ Clinical Resources
▸ Variation
▸ Animal Models
▸ Cellular Pathways

## GEO Profiles

The GEO Profiles database (http://www.ncbi.nlm.nih.gov/geoprofiles) stores individual gene expression profiles from curated DataSets in the Gene Expression Omnibus (GEO) repository (http://www.ncbi.nlm.nih.gov/geo/).

## Worked example 4: GEO Profiles

In this worked example we will look whether there are any expression profiles for the human *BRAF* gene for malignant melanoma.

**STEP 1** – **Go to** the GEO Profiles homepage:
http://www.ncbi.nlm.nih.gov/geoprofiles

STEP 2 – **Enter** 'BRAF[Gene Symbol] AND (human) AND (malignant melanoma)' in the 'Search' text field and **click** [Search]



STEP 3 – **Click on** image for details

Profile    GDS1375 / 206044_s_at / BRAF
Title      Cutaneous malignant melanoma
**Organism** Homo sapiens

**Expression Atlas**

Expression Atlas (http://www.ebi.ac.uk/gxa/home) provides information on gene expression patterns under different biological conditions. It includes both microarray and RNA-seq data. The data is re-analysed to detect interesting expression patterns under the conditions of the original experiment. There are two components to the Expression Atlas, i.e. the Baseline Atlas and the Differential Atlas. The Baseline Atlas displays information about which gene products are present (and at what abundance) in "normal" conditions (e.g. tissue, cell type). The Differential Atlas allows users to identify genes that are up- or down-regulated in a wide variety of different experimental conditions.

**Worked example 5: Expression Atlas**

In this worked example we will have a look what information Expression Atlas contains for the human *BRAF* gene.

**STEP 1** – **Go to** the Expression Atlas homepage:
http://www.ebi.ac.uk/gxa/home



**STEP 2** – **Enter** 'BRAF' in the 'Gene query' text box, **select** 'Organism: Homo sapiens' and **click** [Search]

The resulting page consists of three parts.

The first part contains general information about the *BRAF* gene:



The second part is the Baseline Expression:



And the third part is the Differential Expression:

Detailed information about how to work with Expression Atlas, can be found in the Help section:

## Baseline Atlas at–a–glance



## Differential Atlas at–a–glance

**COSMIC**

Although OMIM is very detailed, it is not comprehensive. COSMIC, the Catalogue of somatic mutations in cancer (http://cancer.sanger.ac.uk/cosmic/), is a specialist resource that aims to have a comprehensive list of genes and their mutations that are involved in cancer. COSMIC curates data from papers in the scientific literature and large scale experimental screens from the Cancer Genome Project (https://www.sanger.ac.uk/research/projects/cancergenome/) at the Sanger Institute. There are several ways to search COSMIC, in the following worked example the most common search interface will be illustrated.

**Worked example 6: COSMIC**

In this worked example we will use COSMIC to list all mutations found in the human *BRAF* gene.

> **STEP 1 – Go to** the COSMIC homepage:
> http://cancer.sanger.ac.uk/cosmic



> **STEP 2 – Enter '**BRAF'
> into the 'Search' text field
> and **click** [SEARCH].

**STEP 5 – Select** 'Mutations' to reveal the molecular details of the mutations

The histogram shows the frequency that the amino acid position has been found to be mutated.

**STEP 6 – Select** 'Tissue'

The 'Mutations' page lists all of the different types of mutations found, including amino acid changes. The frequency of the mutation is shown in the 'Count' column.

Tissue sample summary for BRAF mutations

**DECIPHER**

DECIPHER (DatabasE of Genomic variants and Phenotype in Humans Using Ensembl Resources; https://decipher.sanger.ac.uk/) is a web-based resource and database of genomic variation data from analysis of patient DNA. It documents submicroscopic chromosome abnormalities (microdeletions and duplications) and pathogenic sequence variants (single nucleotide variants - SNVs, Insertions, Deletions, InDels), from over 25,000 patients and maps them to the human genome. In addition it catalogues the clinical characteristics from each patient and maintains a database of microdeletion/duplication syndromes, together with links to relevant scientific reports and support groups.

**Worked example 7: DECIPHER**

In this worked example, we will use DECIPHER to investigate Williams-Beuren Syndrome (WBS), a rare neurodevelopmental disorder (http://en.wikipedia.org/wiki/Williams_syndrome).

> **STEP 1** – **Go to** the DECIPHER homepage:
> https://decipher.sanger.ac.uk/

Syndromes | Gene Disorders

Syndrome List | Karyotype

**STEP 3 – Filter for 'williams'**

Syndromes: 1 to 10 of 70

Filter...

| Syndrome | Location | Interval (Mb) | Grade ? |
|---|---|---|---|
| 1p36 microdeletion syndrome | 1:10001-12840259 | 12.83 | 1 |
| 1q21.1 susceptibility locus for Thrombocytopenia-Absent Radius (TAR) syndrome | 1:145386506-145748067 | 0.36 | 3 |
| Mature Variant Report - RBM8A (c.-21delG) | 1:145507646-145507646 | 0.00 | |
| Mature Variant Report - RBM8A (c.67+32G>C) | 1:145507765-145507765 | 0.00 | |
| 1q21.1 recurrent microdeletion (susceptibility locus for neurodevelopmental disorders) | 1:146533376-147883376 | 1.35 | 3 |
| 1q21.1 recurrent microduplication (possible susceptibility locus for neurodevelopmental disorders) | 1:146533376-147883376 | 1.35 | 3 |
| 2p21 Microdeletion Syndrome | 2:44410451-44589584 | 0.18 | |
| 2p15-16.1 microdeletion syndrome | 2:59285696-61819815 | 2.53 | |
| 2q33.1 deletion syndrome | 2:196925121-205206939 | 8.28 | 1 |
| 2q37 monosomy | 2:239969863-240322643 | 0.35 | 1 |

10

Previous  1  2  3  4  5  6  7  Next

Feedback

Syndromes | Gene Disorders

Syndrome List | Karyotype

**STEP 4 – Click on 'Williams-Beuren Syndrome (WBS)'**

Syndromes: 1 to 1 of 1 (out of 70 total)

williams

| Syndrome | | Interval (Mb) | Grade ? |
|---|---|---|---|
| Williams-Beuren Syndrome (WBS) | 7:72744455-74142672 | 1.40 | 1 |

10

Previous  1  Next

Syndromes » Williams-Beuren Syndrome (WBS)

General information about 'Williams-Beuren Syndrome (WBS)

Overview | Genotype 1 | Citations 9 | Phenotypes 7 | Karyotype

Last modified: 2014-07-02

**Clinical** - Characteristic facial features include periorbital fullness, bulbous nasal tip, long philtrum, wide mouth, full lips, full cheeks and small widely spaced teeth. Individuals have mild to moderate intellectual disability or learning difficulties with relative cognitive strengths in verbal short term memory and in language but extreme weakness in visuospatial construction (writing, drawing, pa[...] include anxiety, attention deficit hyperactivity disorder (ADHD), and overfriendliness. Congenital heart disease occurs in 80%, with t[...]d a smaller proportion having a discrete supravalvular pulmonary stenosis.

The microd[...] which is also mutated in isolated SVAS. Other symptoms include hernias, visual impairment, hypersensitivity to sound, chronic otiti[...]lies, constipation, vomiting, growth deficiency, infantile hypercalcemia, musculoskeletal abnormalities, diabetes and a hoarse voic[...]he distal deletion breakpoint, with hypertension being significantly less prevalent in WBS patients with a deletion that includes NCF[...](p-0.6[...]), a gene coding for the p-47 phox subunit of the NADPH oxidase. This likely arises through life-long reduced angiotensin II-mediated oxidative stress.

**STEP 5 – Click on 'Genotype'**

**Size of deletion** - Three large region-specific LCRs, termed centromeric, medial and telomeric, flank the WBS deletion interval. Each LCR is several hundred kb in length and is comprised of transcriptionally active genes and pseudogenes grouped into discreet blocks known as A, B and C. Most patients (>95%) have a 1.55Mb deletion caused by recombination between centromeric and medial block B copies, which share approximately 99.6% nucleotide identity over many kilobases. There at at hot-spots of recombination: one within a 12 kb region of the GTF2I gene, and one in the distal end of the GTF2IRD2 gene. A few patients (<5%) have a larger deletion (~1.84Mb) caused by recombination between centromeric and medial block A copies.

**Origin of deletion** - Almost one-third (28%) of the transmitting progenitors are heterozygous for an inversion between centromeric and telomeric LCRs which may facilitate the deletion. The deletions are caused by nonhomologous recombination within the LCRs of either the same chromosome 7 (intrachromosomal) or different chromosome 7s (interchromosomal). In each case the chromosomes are envisaged to form loops, thereby allowing the alignment of the two LCRs, the occurrence of recombination, and the excision of the DNA contained within the intervening loop. Approximately 2/3rds of the deletion events are interchromosomal.

**Expert advisors**
Dr. Stephen W. Scherer The Hospital for Sick Children, Toronto, Canada and Dr. Lucy Osborne, University of Toronto, Canada

**Links to support groups:**
www.williams-syndrome.org
www.rarechromo.org

**Links to further information:**
www.ncbi.nlm.nih.gov
www.orpha.net

## Biological ontologies

There is no universal standard terminology in biology and related domains, and term usages may be specific to a species, research area or even a particular research group. Different people may use different terms when referring to the same thing and/or may use the same term for different things. This makes communication and sharing of data more difficult. To try to make everyone using the same term when talking about the same thing, biological ontologies have been developed. Two widely-used biological ontologies are the Gene Ontology (GO) and the Sequence Ontology (SO).

## Gene Ontology (GO)

The Gene Ontology (GO; http://geneontology.org/) consists of three hierarchically structured, controlled vocabularies that describe gene products in terms of their associated biological processes ("What does a gene product do?"), cellular components ("Where does a gene product do what it does?") and molecular functions ("How does a gene product do what it does?") in a species-independent manner.

## Sequence Ontology (SO)

The Sequence Ontology (SO; http://www.sequenceontology.org/) is a set of terms and relationships used to describe the features and attributes of biological sequence.

**Exercises**


**Ensembl Variant Effect Predictor**

Resequencing of the genomic region of the human *CFTR* (cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7) gene (ENSG00000001626) has revealed the following variants (alleles defined in the forward strand):

• Substitution G/A at position 7: 117,530,985
• Substitution T/C at position 7: 117,531,038
• Substitution T/C at position 7: 117,531,068

Use the Ensembl Variant Effect Predictor to answer the following questions:

(a) What genes are affected by these variants?
(b) Do any of the variants result in protein changes?
(c) If so, are these protein changes predicted to be deleterious or damaging?
(d) Have the variants already been annotated in Ensembl?


**OMIM, GEO Profiles and Expression Atlas**

Prostate cancer antigen 3 (*PCA3*, also referred to as *DD3*) is a gene that expresses a non-coding RNA. *PCA3* is only expressed in human prostate tissue, and the gene is highly overexpressed in prostate cancer. Because of its restricted expression profile, the *PCA3* RNA is useful as a tumor marker.

Explore what information is available about *PCA3* in OMIM, GEO Profiles and Expression Atlas. Can you find any evidence that this gene is indeed only expressed in prostate and that it is highly overexpressed in prostate cancer?

**Exercises answers**

**Ensembl Variant Effect Predictor**

(a) The only gene affected by these variants is *CFTR*.

(b) The variant at position 117531038 results in a L/P change at position 138 in two of the CFTR proteins. The variant at position 117531068 results in an I/T change at position 148 in two of the CFTR proteins. The variant at position 117530985 doesn't result in a protein change.

(c) The protein change at position 138 is predicted to be deleterious / damaging. The protein change at position 148 is predicted to be tolerated / benign.

(d) All three variations have been already described and are known as rs1800077, rs1800078 and rs35516286 in dbSNP.

| Uploaded variation | Location | Allele | Gene | Feature | Feature type | Consequence | cDNA position | CDS position | Protein position | Amino acids | Codons | Existing variation | Distance to transcript |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7_117530985_G/A | 7:117530985 | A | ENSG00000001626 | ENST00000446805 | Transcript | downstream_gene_variant | - | - | - | - | - | COSM3632277, COSM3949805, rs1800077 | 7 |
| 7_117530985_G/A | 7:117530985 | A | ENSG00000001626 | ENST00000426809 | Transcript | synonymous_variant | 360 | 360 | 120 | A | GCG/GCA | COSM3632277, COSM3949805, rs1800077 | - |
| 7_117530985_G/A | 7:117530985 | A | ENSG00000001626 | ENST00000003084 | Transcript | synonymous_variant | 492 | 360 | 120 | A | GCG/GCA | COSM3632277, COSM3949805, rs1800077 | - |
| 7_117531038_T/C | 7:117531038 | C | ENSG00000001626 | ENST00000446805 | Transcript | downstream_gene_variant | - | - | - | - | - | rs1800078 | 60 |
| 7_117531038_T/C | 7:117531038 | C | ENSG00000001626 | ENST00000426809 | Transcript | missense_variant | 413 | 413 | 138 | L/P | CTA/CCA | rs1800078 | - |
| 7_117531038_T/C | 7:117531038 | C | ENSG00000001626 | ENST00000003084 | Transcript | missense_variant | 545 | 413 | 138 | L/P | CTA/CCA | rs1800078 | - |
| 7_117531068_T/C | 7:117531068 | C | ENSG00000001626 | ENST00000446805 | Transcript | downstream_gene_variant | - | - | - | - | - | rs35516286, CM962456, CM920145 | 90 |
| 7_117531068_T/C | 7:117531068 | C | ENSG00000001626 | ENST00000426809 | Transcript | missense_variant | 443 | 443 | 148 | I/T | ATT/ACT | rs35516286, CM962456, CM920145 | - |
| 7_117531068_T/C | 7:117531068 | C | ENSG00000001626 | ENST00000003084 | Transcript | missense_variant | 575 | 443 | 148 | I/T | ATT/ACT | rs35516286, CM962456, CM920145 | - |

| Symbol | Symbol source | HGNC ID | Biotype | Transcript support level | SIFT | PolyPhen | GMAF | AFR MAF | EUR MAF | AA MAF | EA MAF | Clinical significance | Somatic status | Pubmed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CFTR | HGNC | HGNC:1884 | protein_coding | 4 | - | - | A:0.0005 | A:0.0020 | - | - | - | - | 1, 1, 0 | - |
| CFTR | HGNC | HGNC:1884 | protein_coding | 5 | - | - | A:0.0005 | A:0.0020 | - | - | - | - | 1, 1, 0 | - |
| CFTR | HGNC | HGNC:1884 | protein_coding | 1 | - | - | A:0.0005 | A:0.0020 | - | - | - | - | 1, 1, 0 | - |
| CFTR | HGNC | HGNC:1884 | protein_coding | 4 | - | - | - | - | - | - | - | - | - | 18716917 |
| CFTR | HGNC | HGNC:1884 | protein_coding | 5 | 0 | 0.962 | - | - | - | - | - | - | - | 18716917 |
| CFTR | HGNC | HGNC:1884 | protein_coding | 1 | 0 | 0.838 | - | - | - | - | - | - | - | 18716917 |
| CFTR | HGNC | HGNC:1884 | protein_coding | 4 | - | - | C:0.0005 | - | C:0.0013 | C:0 | C:0.000814 | not_provided, benign | - | 18716917 |
| CFTR | HGNC | HGNC:1884 | protein_coding | 5 | 0.48 | 0.375 | C:0.0005 | - | C:0.0013 | C:0 | C:0.000814 | not_provided, benign | - | 18716917 |
| CFTR | HGNC | HGNC:1884 | protein_coding | 1 | 0.55 | 0.024 | C:0.0005 | - | C:0.0013 | C:0 | C:0.000814 | not_provided, benign | - | 18716917 |

**OMIM, GEO Profiles and Expression Atlas**

OMIM contains basic information about the *PCA3* gene (official and alternative gene symbols, cytogenetic and genomic location) as well as some

information from a paper that has shown overexpression of this gene in prostate tumors and has suggested that the gene codes for a noncoding RNA.



According to the Baseline Expression in Expression Atlas, the *PCA3* gene is indeed only expressed in the prostate:



GEO Profiles shows a number expression profiles from which it is clear that the *PCA3* gene is highly expressed in prostate cancer. For example:

**Profile**    GDS1439 / 232575_at / PCA3
**Title**      Prostate cancer progression
**Organism** Homo sapiens



GDS1439 / 232575_at / PCA3

- count
- percentile rank within the sample
- value, Detection Call = ABSENT
- rank, Detection Call = ABSENT

237