

BIOL 419/519 Bioinformatics Research – July 16, 2013

Human Genome Annotation – Paul Szauter

I. Introduction to the Personal Genome Project (PGP)

The Personal Genome Project (PGP) is a group of researchers and volunteers who seek to make individual human genomes and health information freely available for public use by anyone. The PGP is an important complement to the biomedical research literature, because the consequences of individual genetic variants can be investigated in the context of an entire genome in a population of individuals who have not been singled out by their disease status.

Go to the PGP website:

<http://www.personalgenomes.org>

Click **PGP Community** in the top navigation menu to go to:

<http://www.personalgenomes.org/community.html>

Click **View public profiles** to go to:

<https://my.personalgenomes.org/users>

Use the selector to **Show 100 entries**. Scroll down to **PGP89** and click **hu011C57**. In the **Complete Genomics** line of the table, click **View report**.

Personal Genome Project Log In ▶

Public data - About -

Public Profile -- hu011C57

Public profile url: <https://my.personalgenomes.org/profile/hu011C57>

Personal Health Records
None added.

Samples

PGP Blood Collection	Sample	Received	Source	Action
	Sample 86486261 (whole blood)	2012-04-26 16:00:00 UTC	Feinstein Institute	Show log
	Sample 60261538 (whole blood)	2012-04-26 16:00:00 UTC	Feinstein Institute	Show log
	Sample 24344335 (whole blood)	2012-05-02 17:43:57 UTC	Coriell	Show log
	Sample 49403627 (whole blood)	2012-05-02 17:43:57 UTC	Coriell	Show log
	Sample 22319525 (whole blood)	2012-05-02 17:43:57 UTC	Coriell	Show log

Uploaded data

Date	Data type	Name	Download	Report
2013-07-15	23andMe	4 files (click to expand)	9.43 GB	
2012-11-09	Complete Genomics	CGI sample GSO1669-DNA_B05 from PGP sample 86486261	Download (234 MB)	View report • male • 2,777,778,141 positions covered • ref. b37
2012-07-20	23andMe	Paul_Szauter_exome_sequence	Download (7.21 MB)	View report • male • 119,413 positions covered • ref. b37
2012-01-26	23andMe	23andMe_SNPs	Download (23.6 MB)	View report • male • 955,463 positions covered • ref. b36

Geographic Information

State: New Mexico
Zip code: 87110

The page that appears displays the **Genome report** by default. These are likely pathogenic and rare variants found in the genome of that individual. Click any allele in the **Variants** column for more information.

hu011C57 - GET-Evidence variant report

[About](#) [Genomes](#) [Guides](#) [Recent changes](#) [Contributors](#) [Download](#)

Variant report for hu011C57

- Data source: CGI sample GS01669-DNA_B05 from PGP sample 86486261
- This report: evidence.personalgenomes.org/genomes/962c2ddfccf53bb761eaf0ce94a334640f5bb436
- Person ID: hu011C57
- public profile: my.personalgenomes.org/profile/hu011C57
- Download: [source data](#) (245 MB), [dbSNP and nsSNP report](#) (118 MB)
- [Show debugging info](#)

Gene search

"GENE" or "GENEA123C":

Log in

OpenID URL:

Genome report				
Insufficiently evaluated variants				
Coverage				
Gene Report				
Metadata				
Show likely pathogenic and rare (<2.5%) pathogenic variants				
Show all				
Show All entries				
Search: <input type="text"/>				
Variant	Clinical Importance	Impact	Allele freq	Summary
RYR2-G1885E	High	Uncertain pathogenic Recessive, Carrier (Heterozygous)	1.8%	Reported to cause arrhythmogenic right ventricular cardiomyopathy when compound heterozygous with G1886S, although this finding is weakened after correcting for multiple hypotheses and it is unclear what penetrance such a genotype might have, if it is causal.
WFS1-C426Y	Moderate	Uncertain pathogenic Dominant, Heterozygous	0.12%	Reported in a single case of familial depression, but no linkage data and no statistical significance.
HFE-C282Y	Low	Well-established pathogenic Recessive, Carrier (Heterozygous)	4.9%	This variant is associated with hereditary haemochromatosis, 80% of patients with that disease are homozygous for this variant. However, the penetrance is low, in Beutler et al. they note that only 1 of their 158 homozygotes met criteria for diagnosis with the condition.
COL4A1-Q1334H	Low	Likely pathogenic Dominant, Heterozygous	32%	This common variant has been associated with arterial stiffness and, in Japanese, a small increased risk of myocardial infarction (MI, a.k.a. heart attack). This last observation supported a dominant effect for this variant and, assuming a lifetime risk of 15% for MI, we estimate carriers have an additional risk of 0.5-3%.
MTRR-I49M	Low	Likely pathogenic Recessive, Homozygous	45%	This common variant (HapMap allele frequency of 31.3%) in a protein involved in folate (B9) and cobalamin (B12) metabolism and is often reported as "MTRR I22M" (an alternative transcript position). Mothers homozygous for this variant are associated with having around a increased chance of a child with Down syndrome (risk of 0.4%, average risk in population is 0.25%). Notably, age plays a far larger role in the rate of Down syndrome (risk is 4.5% for a mother 45-years-of-age), and it is unknown how this variant may combine with the effect of age. There are conflicting reports associating this variant with incidence of neural tube defects, possibly when combined with MTHFR A222V.
KRT5-G138E	Low	Likely pathogenic Unknown, Heterozygous	5.2%	This variant is associated with 1.25x increased risk of basal cell carcinoma (common skin cancer, rarely malignant).
rs5186	Low	Likely pathogenic Unknown, Heterozygous	21%	This common noncoding genetic variant has an allele frequency of ~30% and is associated with an increased risk of hypertension. If ~25% of non-carriers have hypertension, Bonnardeaux et al's data predict ~4% increased risk of hypertension per copy of this variant. This SNP is in the 3' noncoding region of the AGTR1 transcript (angiotensin II type 1 receptor), also known as AT2R1 or AT1R, which is a target of hypertension drugs.
MAD1L1-R59C	Low	Uncertain pathogenic Unknown, Heterozygous	0.36%	Hypothesized to be involved in prostate cancer, but no statistically significant data. Using more detailed variant frequency information, the variant does not appear to be enriched in the cancer samples reported by Tsukasaki et al.

Showing 1 to 8 of 8 entries (filtered from 36 total entries)

The three most important tabs in the top Navigation tabs show you:

- 1. Genome report.** Likely pathogenic and rare variants found in the genome of that individual.
- 2. Insufficiently evaluated variants.** These are variants that need further annotation and evaluation.
- 3. Gene Report.** This shows all genes in which variants have been identified.

In the **Genome Report** table is a **Search** box that lets you search for variants of a specific gene in that individual.

There is a **Gene search box** in the upper right that lets you look for variants of a specific gene among all PGP individuals.

Immediately below the **Gene search box** are **Log in** buttons. If you have a Google or Yahoo email account, you can log in to improve the annotation.

II. Improving annotation for a specific variant

Click the tab for **Insufficiently evaluated variants**. In the Search box in the Genome report table, type **CFTR**. A single variant appears. Click the allele.

Genome report Insufficiently evaluated variants Coverage Gene Report Metadata

Show 100 entries Search: CFTR

Variant	Prioritization score	Allele freq	Num of articles	Zygotity and Prioritization Score Reasons
CFTR-R75Q	3	?		Heterozygous. Polyphen 2: Unknown, Testable gene in GeneTests with associated GeneReview

Showing 1 to 1 of 1 entries (filtered from 3,221 total entries)

The allele is CFTR-R75Q. This is a missense allele in which the 75th amino acid of the CFTR protein is changed from Arginine (R) to Glutamine (Q).

To evaluate whether this allele is likely to be harmful, we use PolyPhen-2. Open a new browser window to the PolyPhen-2 site from the Tools page on the course website:

<http://genetics.bwh.harvard.edu/pph2/>

Enter **CFTR** as the Protein, **75** as the position, and select **R** as AA1 and **Q** as AA2. Click **Submit Query**.

Query Data	
Protein or SNP identifier	<input type="text" value="CFTR"/>
Protein sequence in FASTA format	<input type="text"/>
Position	<input type="text" value="75"/>
Substitution	AA ₁ A R N D C E Q G H I L K M F P S T W Y V AA ₂ A R N D C E Q G H I L K M F P S T W Y V
Query description	<input type="text"/>

[Display advanced query options](#)

When the job is done running, click the link to go to the screen shown below. Click the plus sign to expand the **Multiple sequence alignment**.

PolyPhen-2 report for P13569 R75Q (rs1800076)

Query

Protein Acc	Position	AA ₁	AA ₂	Description
P13569	75	R	Q	Canonical; RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel conductance-controlling ATPase; EC=3.6.3.49; AltName: Full=cAMP-dependent chloride channel; Length: 1480

Results

Prediction/Confidence *PolyPhen-2 v2.2.2r398*

HumDiv

This mutation is predicted to be **PROBABLY DAMAGING** with a score of **1.000** (sensitivity: **0.00**; specificity: **1.00**)

0.00 0.20 0.40 0.60 0.80 1.00

HumVar

Details

Multiple sequence alignment *UniProtKB/UniRef100 Release 2011_12 (14-Dec-2011)*

```

QUERY          VDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q2QL83#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp F6RU09#1     ADSADNLSEKLE-RE-WDRELV---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q2QLH0#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q00554#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q6PQ22#1     SDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q07E16#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp F1SJE0#1     SDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q00PJ2#1     TDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q07E42#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp F1SJE1#1     SDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp A0M8T4#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q108U0#1     VDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp G5B3Q7#1     ADSADNLSEKLE-RE-WDRELV---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp Q07E05#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp G1UHC1#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp G3SWU6#1     VDSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
sp G1SZW4#1     ADSADNLSEKLE-RE-WDREIA---SKKNEKLINALR R CFEWRRFMFYGIILLYLGEV--TKAVCPDLLLGRIIASYPD
  
```

Shown are 75 amino acids surrounding the mutation position (marked with a black box). An interactive version of the complete alignment is [also available](#).

This allele is rated as **Probably damaging**. Each line in the multiple alignment represents a different species. Click any line to see the entry for that species.

We can use OMIM to see if this variant has already been described. Open a new window to OMIM using the link from the Tools page of the course website:

<http://www.omim.org>

Enter **CFTR** in the search box and click **Search**.

The top link takes you to the gene entry for CFTR. At the top of the page are links to specific disease entries associated with the CFTR gene.

Scroll down to the **Allelic Variants** section or just search for **ARG75GLN**. This search fails, so try links in the **Variation** menu in the right margin. The **Locus Specific DBs** link works; search for ARG75GLN in the database there.

There is a citation listed, which you can locate through PubMed (PMID: 1710599). This paper is cited by a more recent paper (PMID: 20977904). You can use **summarize the information in this paper** on the page for the CFTR-R75Q allele (click the evaluating evidence link for more information).

III. Searching for all variants of a specific gene in the PGP

In the **Gene Search box**, enter **CFTR**.

This returns a multiple-page table of CFTR variants. The **Genomes** column contains links to individual genomes with that variant. Note that some of these alleles have summaries written by annotators.

IV. Approaches to improving annotation at the PGP

There are several approaches to improving annotation at the PGP, outlined below.

A. Genomic Checkup. The American College of Medical Genetics (ACMG) has released a list of genes to be investigated as incidental findings when a person's genome is analyzed. Use this gene list to perform a checkup on an individual in the PGP. Do you have any reportable findings?

B. Newborn Screening. Newborns are screened for treatable metabolic disorders. The genes associated with these disorders are known. You can screen individuals in the PGP to see if they carry variant alleles of any of these genes.

C. Insufficiently Evaluated Variants. For any individual, look at insufficiently evaluated variants. Pick variants whose frequency is known, but below 1% (or 5% if you are ambitious).

D. Literature Based. Use OMIM or the primary literature to find genes that are interesting. Look for variants among individuals in the PGP.

V. Report Interesting findings

Keep track of the genes and variants that you have evaluated, even those that you have found uninteresting. If you find something interesting, share it with the class. If you are not confident about annotating the variant yourself, collect and submit your information to pzauter@unm.edu for further evaluation.