



1390 Shorebird Way
Mountain View, CA 94043
www.23andme.com

Exome Results & Raw Data Summary

Generated on: 4/26/2012

Congratulations! Your exome has been sequenced and your data is ready for you to download. We have also included this overview of your data to get you started on your exome exploration. Here are a few important points about your exome data:

- Two types of files are available for download: 1) the aligned sequencing reads in BAM format, 2) a file containing variant calls (VCF file).
- The raw data VCF file is a preliminary draft of your exome. Our ability to call variants, especially indels, is greatly improved with each additional exome added to our database. Moreover we will build upon this protocol to include additional steps such as custom treatment of the sex chromosomes. To this end we will update your VCF file at the end of the pilot. We will contact you when this data is available.

Your exome at a glance:

[Your exome in numbers](#)

[Characterizing your variants](#)

[How rare are your variants?](#)

[Filtering your variants](#)

[See selected variants](#)

[Appendix](#)

The Exome Service is a pilot project, and this report contains preliminary data only. 23andMe does not represent that all of this information is accurate. **In this report we have used 1000 Genome Project data to report frequencies of variants to determine how common or rare a particular variant is.** We have also only provided information about a subset of the many gene-disrupting variants present in the human genome, in a chosen set of genes. Sequencing was performed such that the total number of bases read was at least 80X the size of the exome. As described in the Exome Terms of Use, 23andMe will not be providing the reports and explanations that 23andMe typically provides to customers with respect to their genotyping results for this data. 23andMe Services are for research, informational, and educational use only. We do not provide medical advice. Please keep in mind that genetic information you share with others could be used against your interests.

Your exome in numbers

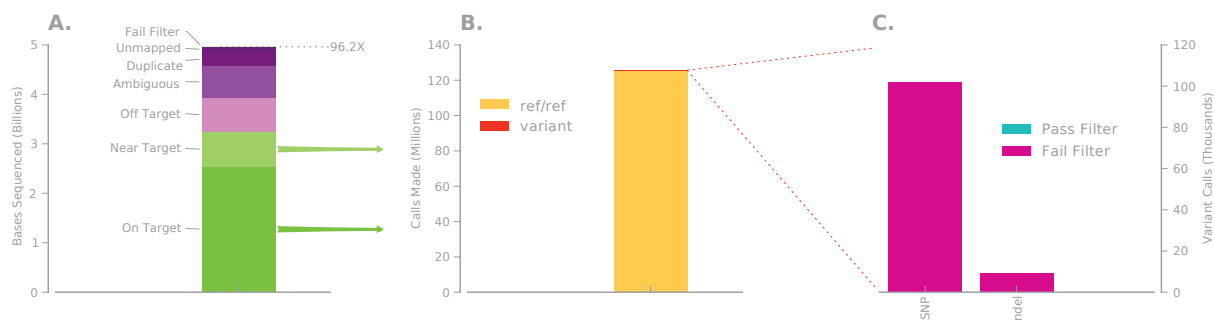


Figure 1: Getting from raw reads to called variants. A) The number of bases obtained by sequencing your exome. The top line indicates total coverage. B) Total number of called bases in your exome. The vast majority are the same as the reference genome. C) An expansion of the small sliver of variants depicted in B. These are the variants present in your VCF file.

Welcome to your exome. Your exome is the 50 million DNA bases of your genome containing the information necessary to encode all your proteins. Your exome data consists of two parts, the raw data (both aligned and unaligned Illumina reads, fig1A) and a draft of the variants present in your exome (fig1C). While this draft is provisional and we will be improving upon it, we wanted to allow you to dig in to your exome as soon as possible so you can tell us what you think is important and should be included.

To create the first draft of your exome we implemented the Broad Institute's "Best Practice" protocol for exome sequencing analysis. You can read a detailed description of it [here](#) (for brief summary see [Appendix](#)).

Characterizing your variants

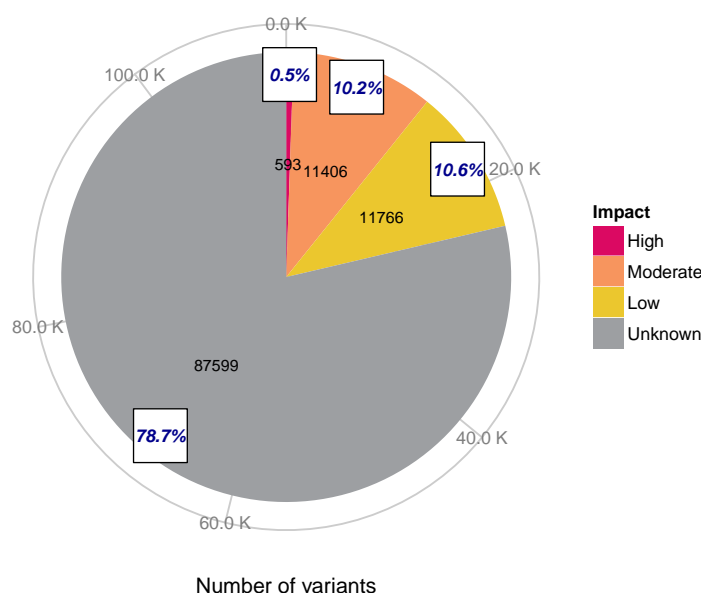


Figure 2: Predicting impact of variants on gene function. An overview of your variants and their predicted impact on gene function.

The variants in your VCF file are the positions in your genome that differ from the reference genome. Most of these variants are likely to be functionally neutral and unlikely to cause any severe disorders. Pinpointing genuine disease mutations is still challenging and we used a number of software tools to identify those that may be functionally important. We estimated the impact a variant has on gene function based on the severity of its effect on the gene product:

High impact:

Frame shift Insertion or deletion of bases, not multiple of 3.

Splice site Variant at the 'splicing site' may disrupt the consensus splicing site sequence.

Stop gain Premature termination of peptides, which would disable protein function.

Start loss Loss of the start codon.

Stop loss Loss of the stop codon.

Moderate impact:

Nonsynonymous substitution Non-conservative change altering an amino acid in a protein.

Codon insertion or deletion Insertion or deletion of bases, multiple of 3.

Low impact:

Synonymous substitution Variant that does not alter the amino acid sequence due to codon degeneracy.

Start gain Variant resulting in the gain of a start codon.

Synonymous stop Variant changing one stop codon into another.

Unknown impact: Variants unlikely to affect gene products.

How rare are your variants?

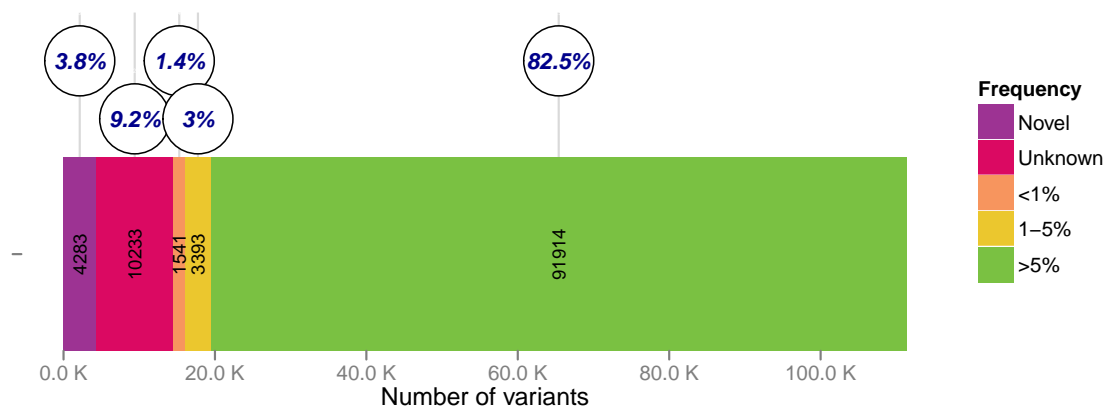


Figure 3: Variant frequencies. The allele frequencies of the variants in your exome. Unknown: allele is present in a public database but no frequency data was available.

One of the advantages of exome sequencing is that we can detect sequence variants that are unique to you! By comparing your variants to all those that have been discovered so far, we can divide your variants into the following categories:

- **novel** variant hasn't been observed in current public sequence databases
- **unknown** variant has been observed in public databases but allelic frequency has not been calculated and therefore is not available
- **rare** variant with allelic frequency <1%
- **somewhat rare** variant with frequency 1-5%
- **common** frequency of the variant is greater than 5%

One of the most comprehensive human variation public datasets is maintained by the 1000 Genomes Project. We use 1000 Genomes Project data (project release: 08-26-2011) to report frequencies of alleles found in your exome, including reporting if it is absent from the public database (*i.e.* a novel variant).

Filtering your variants

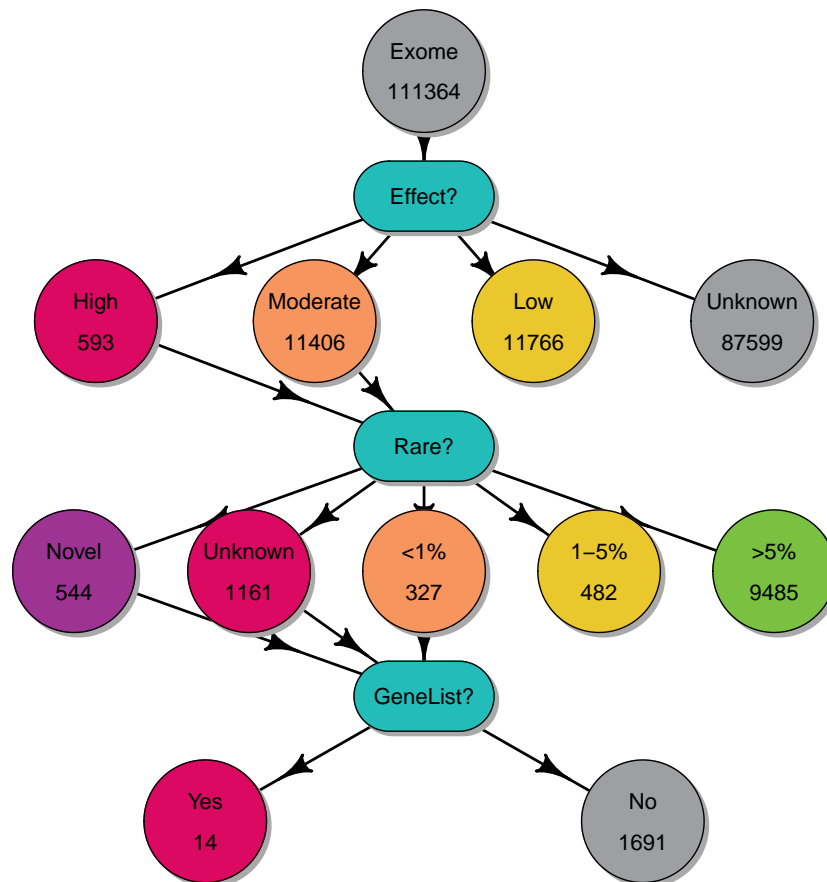


Figure 4: Variant filtering decision tree. A graphical representation of the filtering process that was used to generate your short list of variants of interest.

Most sequence variants in your exome are likely to be neutral and do not cause any severe disorders. A filtering process is often undertaken to prioritize variants discovered through sequencing. To identify potentially interesting and relevant variants with potential functional effects (contributing to disease and other phenotypes of interest) we used three consecutive filters, depicted in the figure above: (1) effect of the variant on the gene product; (2) allele frequency of the variant; (3) location of the variant in one of 592 genes involved in Mendelian disorders (at this point we also exclude indels and variants on the sex chromosomes).

We hope you find this initial list of variants interesting and that it will help you in your journey through your exome. This short list of variants only scratches the surface of what your genome contains and is just the beginning of where your data can take you. Have fun!

List of selected variants

Variant 1:	Gene: DHCR7 Your genotype: C/T Location: chr11:71155265		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00880	dbSNP: rs140748737	
Quality:	Genotype quality: 59.73	Coverage depth: 6	
Details:	Gene description: 7-dehydrocholesterol reductase Transcript: ENST00000529990 AA change: R12H EntrezId: 1717 EnsemblId: ENSG00000172893 UniProt: Q9UBM7 OMIM: 602858		

NON SYNONYMOUS CODING (R12H)



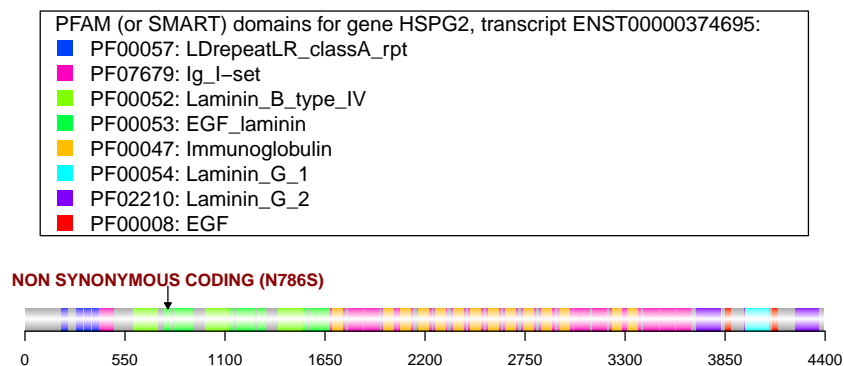
Variant 2:	Gene: LRPPRC Your genotype: C / T Location: chr2:44116923		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00320	dbSNP: rs147302249	
Quality:	Genotype quality: 99	Coverage depth: 79	
Details:	Gene description: leucine-rich pentatricopeptide repeat containing Transcript: ENST00000260665 AA change: A1360T EntrezId: 10128 EnsemblId: ENSG00000138095 UniProt: P42704 OMIM: 607544		

PFAM (or SMART) domains for gene LRPPRC, transcript ENST00000260665:
■ PF01535: Pentatricopeptide_repeat

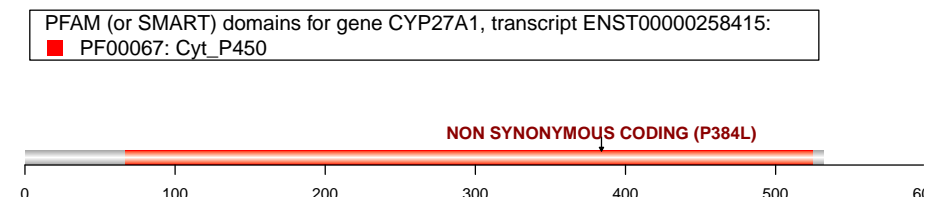
NON SYNONYMOUS CODING (A1360T)



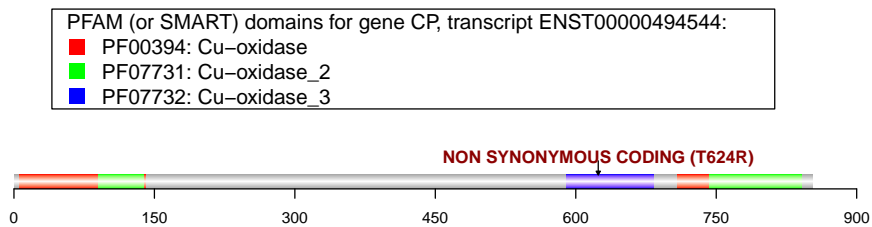
Variant 3:	Gene: HSPG2 Your genotype: T/C Location: chr1:22205601		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00190	dbSNP: rs143736974	
Quality:	Genotype quality: 99	Coverage depth: 48	
Details:	Gene description: heparan sulfate proteoglycan 2 Transcript: ENST00000374695 AA change: N786S EntrezId: 3339 EnsemblId: ENSG00000142798 UniProt: P98160 OMIM: 142461		



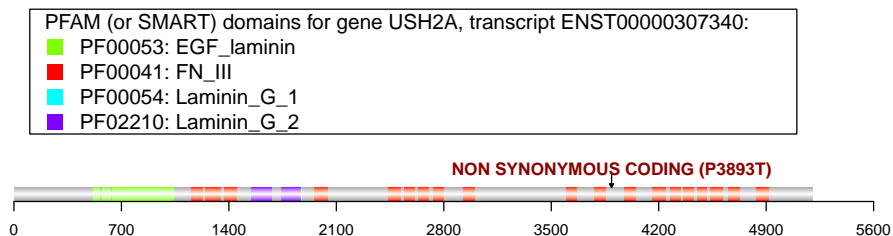
Variant 4:	Gene: CYP27A1 Your genotype: C/T Location: chr2:219678877		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00820	dbSNP: rs41272687	
Quality:	Genotype quality: 99	Coverage depth: 35	
Details:	Gene description: cytochrome P450, family 27, subfamily A, polypeptide 1 Transcript: ENST00000258415 AA change: P384L EntrezId: 1593 EnsemblId: ENSG00000135929 UniProt: Q02318 OMIM: 606530		



Variant 5:	Gene: CP Your genotype: G/C Location: chr3:148899824		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00460		dbSNP: rs56033670
Quality:	Genotype quality: 99		Coverage depth: 39
Details:	Gene description: ceruloplasmin (ferroxidase) Transcript: ENST00000494544 AA change: T624R EntrezId: 1356 EnsemblId: ENSG00000047457 UniProt: P00450 OMIM: 117700		

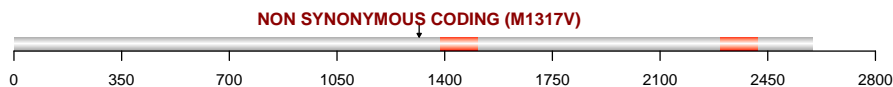


Variant 6:	Gene: USH2A Your genotype: G / T Location: chr1:215914751		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00780		dbSNP: rs41303285
Quality:	Genotype quality: 99		Coverage depth: 76
Details:	Gene description: Usher syndrome 2A (autosomal recessive, mild) Transcript: ENST00000307340 AA change: P3893T EntrezId: 7399 EnsemblId: ENSG00000042781 UniProt: O75445 OMIM: 608400		

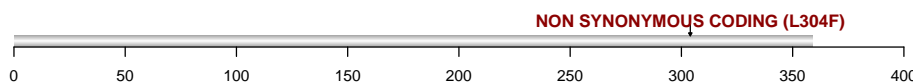


Variant 7:	Gene: ABCA12 Your genotype: T/C Location: chr2:215852398		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00000		dbSNP: rs145178648
Quality:	Genotype quality: 99		Coverage depth: 18
Details:	Gene description: ATP-binding cassette, sub-family A (ABC1), member 12 Transcript: ENST00000272895 AA change: M1317V EntrezId: 26154 EnsemblId: ENSG00000144452 UniProt: Q86UK0 OMIM: 607800		

PFAM (or SMART) domains for gene ABCA12, transcript ENST00000272895:
■ PF00005: ABC_transporter-like



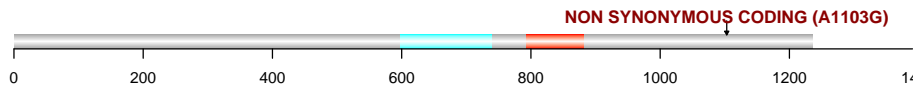
Variant 8:	Gene: CRTAP Your genotype: C/T Location: chr3:33174163		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00550		dbSNP: rs115198029
Quality:	Genotype quality: 99		Coverage depth: 90
Details:	Gene description: cartilage associated protein Transcript: ENST00000449224 AA change: L304F EntrezId: 10491 EnsemblId: ENSG00000170275 UniProt: O75718 OMIM: 605497		



Variant 9:	Gene: RPGRIP1L Your genotype: G/C Location: chr16:53653005		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00640		dbSNP: rs139974543
Quality:	Genotype quality: 99		Coverage depth: 98
Details:	<div>Gene description: RPGRIP1-like</div> <div>Transcript: ENST00000262135</div> <div>EntrezId: 23322</div> <div>UniProt: Q68CZ1</div> <div>AA change: A1103G</div> <div>EnsemblId: ENSG00000103494</div> <div>OMIM: 610937</div>		

PFAM (or SMART) domains for gene RPGRIP1L, transcript ENST00000262135:

- PF11618: DUF3250
- PF00168: C2_Ca-dep



Variant 10:	Gene: ANK2 Your genotype: T/C Location: chr4:114279628		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00320	dbSNP: rs36210417	
Quality:	Genotype quality: 99	Coverage depth: 77	
Details:	Gene description: ankyrin 2, neuronal Transcript: ENST00000505342 EntrezId: 287 UniProt: Q01484		AA change: I295T EnsemblId: ENSG00000145362 OMIM: 106410

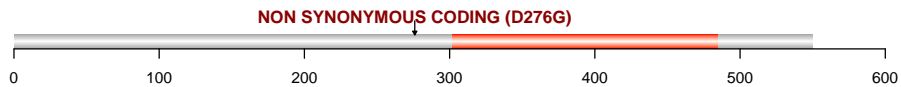
PFAM (or SMART) domains for gene ANK2, transcript ENST00000505342:

- PF00531: Death



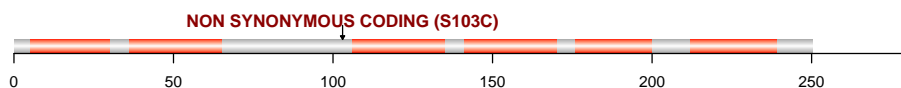
Variant 11:	Gene: MKS1 Your genotype: T/C Location: chr17:56290344		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 5e-04	dbSNP: rs151023718	
Quality:	Genotype quality: 99	Coverage depth: 32	
Details:	Gene description: Meckel syndrome, type 1 Transcript: ENST00000537529 AA change: D276G EntrezId: 54903 EnsemblId: ENSG00000011143 UniProt: Q9NXB0 OMIM: 609883		

PFAM (or SMART) domains for gene MKS1, transcript ENST00000537529:
■ PF07162: B9



Variant 12:	Gene: NEB Your genotype: G/C Location: chr2:152394444		
Effect:	Impact: NON SYNONYMOUS CODING	Type: MODERATE	
Frequency:	1KGenomes: 0.00620		dbSNP: rs62167164
Quality:	Genotype quality: 99		Coverage depth: 66
Details:	Gene description: nebulin Transcript: ENST00000420924 EntrezId: 4703 UniProt: P20929		AA change: S103C EnsemblId: ENSG00000183091 OMIM: 161650

PFAM (or SMART) domains for gene NEB, transcript ENST00000420924:
■ PF00880: Nebulin_35r-motif



Variant 13: Gene: [LRP2](#) Your genotype: [A/G](#) Location: chr2:169989127

Effect: **Impact:** NON SYNONYMOUS CODING **Type:** MODERATE

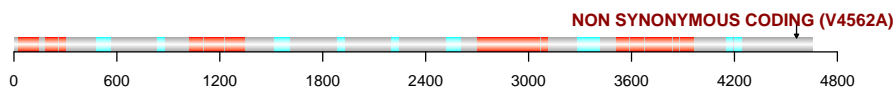
Frequency: **1KGenomes:** 5e-04 **dbSNP:** [rs142245618](#)

Quality: **Genotype quality:** 99 **Coverage depth:** 53

Details: **Gene description:** low density lipoprotein receptor-related protein 2
Transcript: [ENST00000263816](#) **AA change:** V4562A
EntrezId: 4036 **EnsemblId:** [ENSG00000081479](#)
UniProt: [P98164](#) **OMIM:** [600073](#)

PFAM (or SMART) domains for gene LRP2, transcript ENST00000263816:

- PF00057: LDrepeatLR_classA_rpt
- PF00058: LDLR_classB_rpt



Variant 14: Gene: [GJB2](#) Your genotype: [A/G](#) Location: chr13:20763620

Effect: **Impact:** NON SYNONYMOUS CODING **Type:** MODERATE

Frequency: **1KGenomes:** 0.00960 **dbSNP:** [rs35887622](#)

Quality: **Genotype quality:** 99 **Coverage depth:** 65

Details: **Gene description:** gap junction protein, beta 2, 26kDa
Transcript: [ENST00000382844](#) **AA change:** M34T
EntrezId: 2706 **EnsemblId:** [ENSG00000165474](#)
UniProt: [P29033](#) **OMIM:** [121011](#)

PFAM (or SMART) domains for gene GJB2, transcript ENST00000382844:

- PF00029: Connexin_N
- PF10582: Connexin_CCC



Appendix

To create the first draft of your exome we implemented the Broad Institute's "Best Practice" protocol for exome sequencing analysis. You can read a detailed description of it [here](#), however a brief summary of it follows:

1. We took your raw reads and aligned them against the reference genome (these are the alignments available in the BAM file of the encrypted download).
2. We used these alignments to identify probable contamination (unaligned reads) and artifacts of sample preparation (PCR duplicates) which are then removed from subsequent steps.
3. From this point on we focus on the reads that align either to one of the exons or within the regions 250 bases up and downstream of it.
4. To improve the quality of the alignments we carry out a more accurate alignment of the reads that overlap known indels or are likely to contain indels themselves.
5. We also recalibrate the base quality scores of the reads to bring them in line with the empirically-determined values.
6. Using these realigned+recalibrated reads we generate allele calls at every position with enough high-quality data and filter out those that are homozygous for the allele present in the reference genome (the vast majority of these are at such a high frequency in the population they're unlikely to be interesting). The remaining SNP and indel calls (variants) are the ones available in the VCF file that you downloaded.
7. As yet no sequencing technology is 100% accurate and the highly duplicated nature of the human genome makes variant calling a challenging task. Consequently, a small proportion of the variant calls in your VCF are likely to be incorrect. To reduce this proportion we applied the filters recommended by the Broad Institute to remove technical artifacts. Variants that pass all filters are marked in your VCF file with a PASS. As the exome pilot progresses and we gather more data we will be able to use more advanced techniques identify potential errors and improve the quality of your exome.