



# Exome Results & Raw Data Summary

Generated on: June 20, 2012

1390 Shorebird Way  
Mountain View, CA 94043  
[www.23andme.com](http://www.23andme.com)

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?](#)

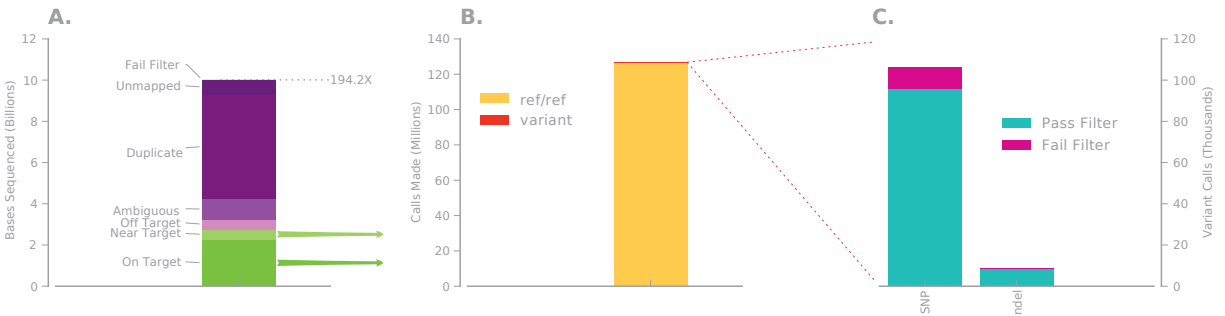
[Filtering your variants](#)

[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

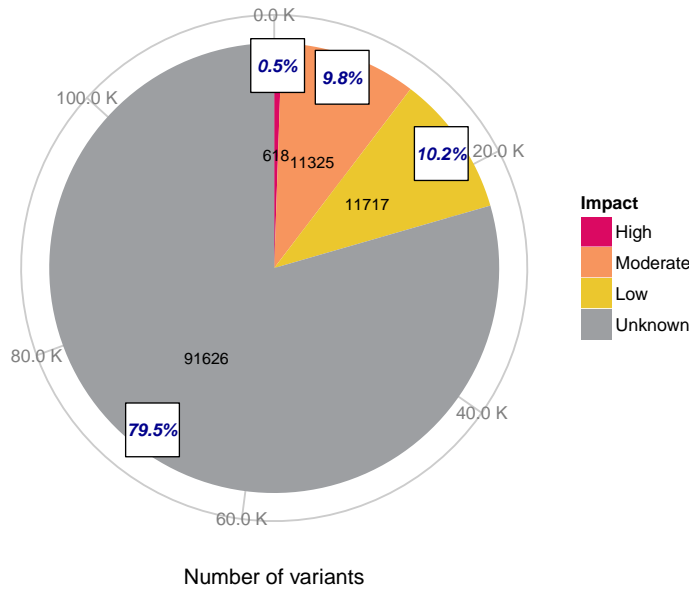


**Figure 1: Getting from raw reads to called variants.** A) The number of bases obtained by sequencing your exome. The top line indicates total coverage. B) Total number of called bases in your exome. The vast majority are the same as the reference genome. C) An expansion of the small sliver of variants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

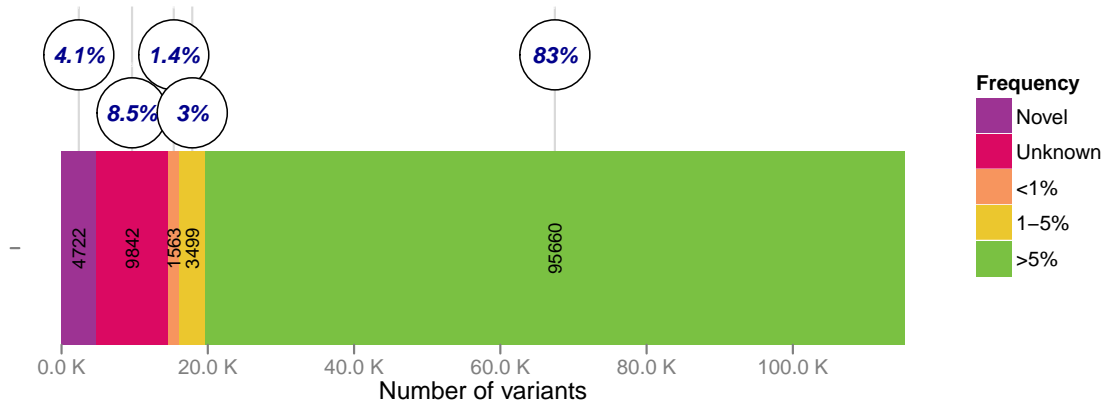
**Synonymous substitution** Variant that does not alter the amino acid sequence due to codon degeneracy.

**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?



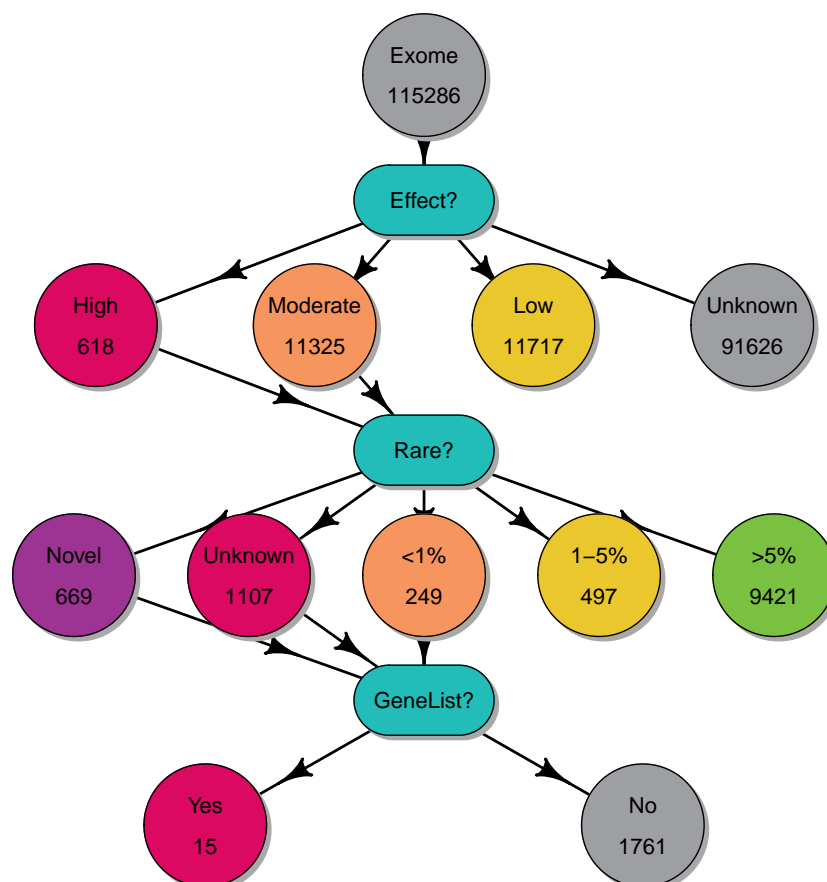
**Figure 3: Variant frequencies.** The allele frequencies of the variants in your exome. Unknown: allele is present in a public database but no frequency data was available.

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.

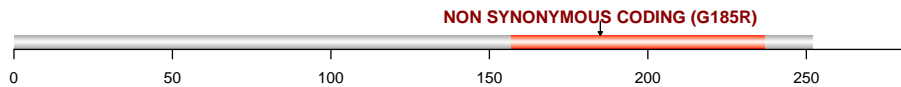
Most sequence variants in your exome are likely to be neutral and do not cause any severe disorders. A filtering process is often undertaken to prioritize variants discovered through sequencing. To identify potentially interesting and relevant variants with potential functional effects (contributing 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

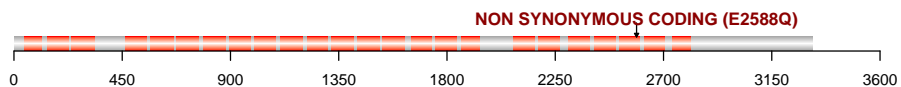
<b>Variant 1:</b>	<b>Gene:</b> <a href="#">SP110</a> <b>Your genotype:</b> C/T <b>Location:</b> chr2:231042873
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00690 <b>dbSNP:</b> <a href="#">rs149485401</a>
<b>Quality:</b>	<b>Genotype quality:</b> 87.8 <b>Coverage depth:</b> 10
<b>Details:</b>	<b>Gene description:</b> SP110 nuclear body protein <b>Transcript:</b> <a href="#">ENST00000338556</a> <b>AA change:</b> G185R <b>EntrezId:</b> 3431 <b>EnsemblId:</b> <a href="#">ENSG00000135899</a> <b>UniProt:</b> <a href="#">Q9HB58</a> <b>OMIM:</b> <a href="#">604457</a>

PFAM (or SMART) domains for gene SP110, transcript ENST00000338556:  
■ PF01342: SAND\_dom



<b>Variant 2:</b>	<b>Gene:</b> <a href="#">CDH23</a> <b>Your genotype:</b> G/C <b>Location:</b> chr10:73563067
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00670 <b>dbSNP:</b> <a href="#">rs41281338</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 26
<b>Details:</b>	<b>Gene description:</b> cadherin-related 23 <b>Transcript:</b> <a href="#">ENST00000398855</a> <b>AA change:</b> E2588Q <b>EntrezId:</b> 64072 <b>EnsemblId:</b> <a href="#">ENSG00000107736</a> <b>UniProt:</b> <a href="#">Q9H251</a> <b>OMIM:</b> <a href="#">605516</a>

PFAM (or SMART) domains for gene CDH23, transcript ENST00000398855:  
■ PF00028: Cadherin

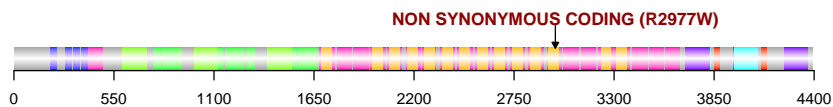


<b>Variant 3:</b>	<b>Gene:</b> <a href="#">EVC</a> <b>Your genotype:</b> <a href="#">G/A</a> <b>Location:</b> chr4:5733317
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 9e-04 <b>dbSNP:</b> <a href="#">rs41269549</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 94
<b>Details:</b>	<b>Gene description:</b> Ellis van Creveld syndrome <b>Transcript:</b> <a href="#">ENST00000264956</a> <b>AA change:</b> D184N <b>EntrezId:</b> 2121 <b>EnsemblId:</b> <a href="#">ENSG00000072840</a> <b>UniProt:</b> <a href="#">P57679</a> <b>OMIM:</b> <a href="#">604831</a>



<b>Variant 4:</b>	<b>Gene:</b> <a href="#">HSPG2</a> <b>Your genotype:</b> <a href="#">G/A</a> <b>Location:</b> chr1:22168855
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00340 <b>dbSNP:</b> <a href="#">rs114851469</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 16
<b>Details:</b>	<b>Gene description:</b> heparan sulfate proteoglycan 2 <b>Transcript:</b> <a href="#">ENST00000374695</a> <b>AA change:</b> R2977W <b>EntrezId:</b> 3339 <b>EnsemblId:</b> <a href="#">ENSG00000142798</a> <b>UniProt:</b> <a href="#">P98160</a> <b>OMIM:</b> <a href="#">142461</a>

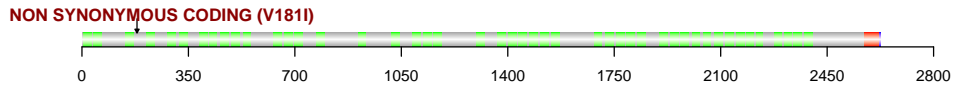
- PFAM (or SMART) domains for gene HSPG2, transcript ENST00000374695:
- PF00057: LDrepeatLR\_classA\_rpt
  - PF07679: Ig\_I-set
  - PF00052: Laminin\_B\_type\_IV
  - PF00053: EGF\_laminin
  - PF00047: Immunoglobulin
  - PF00054: Laminin\_G\_1
  - PF02210: Laminin\_G\_2
  - PF00008: EGF



<b>Variant 5:</b>	<b>Gene:</b> <a href="#">NEB</a> <b>Your genotype:</b> C/T <b>Location:</b> chr2:152425820
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00100 <b>dbSNP:</b> <a href="#">rs149881695</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 66
<b>Details:</b>	<b>Gene description:</b> nebulin <b>Transcript:</b> <a href="#">ENST00000397342</a> <b>AA change:</b> V181I <b>EntrezId:</b> 4703 <b>EnsemblId:</b> <a href="#">ENSG00000183091</a> <b>UniProt:</b> <a href="#">P20929</a> <b>OMIM:</b> <a href="#">161650</a>

PFAM (or SMART) domains for gene NEB, transcript ENST00000397342:

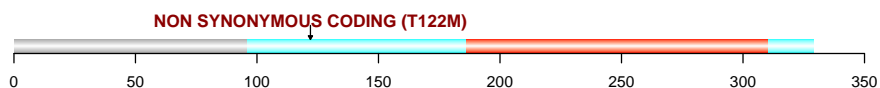
- PF00880: Nebulin\_35r-motif
- PF07653: SH3\_2
- PF00018: SH3\_domain



<b>Variant 6:</b>	<b>Gene:</b> <a href="#">BCKDHA</a> <b>Your genotype:</b> C/T <b>Location:</b> chr19:41920030
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00560 <b>dbSNP:</b> <a href="#">rs34442879</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 19
<b>Details:</b>	<b>Gene description:</b> branched chain keto acid dehydrogenase E1, alpha polypeptide <b>Transcript:</b> <a href="#">ENST00000542943</a> <b>AA change:</b> T122M <b>EntrezId:</b> 593 <b>EnsemblId:</b> <a href="#">ENSG00000248098</a> <b>UniProt:</b> <a href="#">P12694</a> <b>OMIM:</b> <a href="#">608348</a>

PFAM (or SMART) domains for gene BCKDHA, transcript ENST00000542943:

- PF00676: DH\_E1
- PF00456: Transketolase\_N





**Variant 7:** Gene: [GLE1](#) Your genotype: **G/A** Location: chr9:131287573

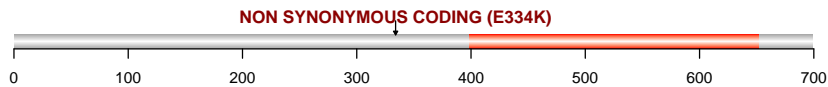
**Effect:** Impact: NON SYNONYMOUS CODING Type: MODERATE

**Frequency:** 1KGenomes: 0.00460 dbSNP: [rs138310419](#)

**Quality:** Genotype quality: 99 Coverage depth: 14

**Details:** Gene description: GLE1 RNA export mediator homolog (yeast)  
Transcript: [ENST00000309971](#) AA change: E334K  
EntrezId: 2733 EnsemblId: [ENSG00000119392](#)  
UniProt: [Q53GS7](#) OMIM: 603371

PFAM (or SMART) domains for gene GLE1, transcript ENST00000309971:  
■ PF07817: GLE1



**Variant 8:** Gene: [MSH2](#) Your genotype: **G/A** Location: chr2:47643457

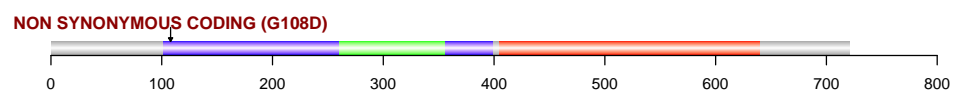
**Effect:** Impact: NON SYNONYMOUS CODING Type: MODERATE

**Frequency:** 1KGenomes: 0.00910 dbSNP: [rs4987188](#)

**Quality:** Genotype quality: 99 Coverage depth: 31

**Details:** Gene description: mutS homolog 2, colon cancer, nonpolyposis type 1 (E. coli)  
Transcript: [ENST00000413880](#) AA change: G108D  
EntrezId: 4436 EnsemblId: [ENSG00000095002](#)  
UniProt: [P43246](#) OMIM: 609309

PFAM (or SMART) domains for gene MSH2, transcript ENST00000413880:  
■ PF05192: DNA\_mismatch\_repair\_MutS\_core  
■ PF05190: DNA\_mismatch\_repair\_MutS\_clamp  
■ PF00488: DNA\_mismatch\_repair\_MutS\_C



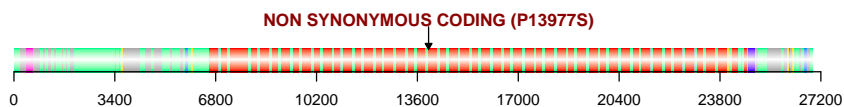
<b>Variant 9:</b>	<b>Gene:</b> <a href="#">PTCH1</a> <b>Your genotype:</b> C/T <b>Location:</b> chr9:98212185
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 9e-04 <b>dbSNP:</b> <a href="#">rs113663584</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 41
<b>Details:</b>	<b>Gene description:</b> patched 1 <b>Transcript:</b> <a href="#">ENST00000418258</a> <b>AA change:</b> G1012S <b>EntrezId:</b> 5727 <b>EnsemblId:</b> <a href="#">ENSG00000185920</a> <b>UniProt:</b> <a href="#">Q13635</a> <b>OMIM:</b> <a href="#">601309</a>

PFAM (or SMART) domains for gene PTCH1, transcript ENST00000418258:  
■ PF02460: Patched  
■ PF12349:



<b>Variant 10:</b>	<b>Gene:</b> <a href="#">TTN</a> <b>Your genotype:</b> G/A <b>Location:</b> chr2:179441932
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING <b>Type:</b> MODERATE
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00240 <b>dbSNP:</b> <a href="#">rs55980498</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99 <b>Coverage depth:</b> 112
<b>Details:</b>	<b>Gene description:</b> titin <b>Transcript:</b> <a href="#">ENST00000356127</a> <b>AA change:</b> P13977S <b>EntrezId:</b> 7273 <b>EnsemblId:</b> <a href="#">ENSG00000155657</a> <b>UniProt:</b> <a href="#">Q8WZ42</a> <b>OMIM:</b> <a href="#">188840</a>

PFAM (or SMART) domains for gene TTN, transcript ENST00000356127:  
■ PF07679: Ig\_I-set  
■ PF09042: Titin\_Z  
■ PF00047: Immunoglobulin  
■ PF07686: Ig\_V-set  
■ PF00041: FN\_III  
■ PF00069: Se/Thr\_kinase-like\_dom  
■ PF07714: Ser-Thr/Tyr\_kinase



Variant 11: Gene: [ITGB4](#) Your genotype: C/T Location: chr17:73738809

Effect: Impact: NON SYNONYMOUS CODING Type: MODERATE

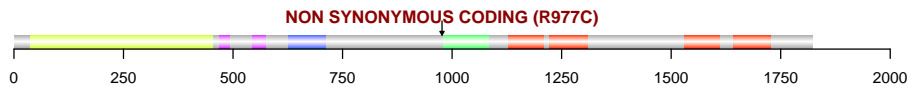
Frequency: 1KGenomes: 0.00140 dbSNP: [rs145976111](#)

Quality: Genotype quality: 99 Coverage depth: 27

Details: Gene description: integrin, beta 4  
Transcript: [ENST00000200181](#) AA change: R977C  
EntrezId: 3691 EnsemblId: [ENSG00000132470](#)  
UniProt: [P16144](#) OMIM: [147557](#)

PFAM (or SMART) domains for gene ITGB4, transcript ENST00000200181:

- PF00362: Integrin\_bsu\_N
- PF07974: EGF\_extracell
- PF07965: Integrin\_bsu\_tail
- PF03160: Calx\_beta
- PF00041: FN\_III



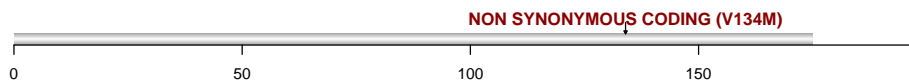
Variant 12: Gene: [MLH1](#) Your genotype: G/A Location: chr3:37092019

Effect: Impact: NON SYNONYMOUS CODING Type: MODERATE

Frequency: 1KGenomes: 4e-04 dbSNP: [rs35831931](#)

Quality: Genotype quality: 99 Coverage depth: 27

Details: Gene description: mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)  
Transcript: [ENST00000421440](#) AA change: V134M  
EntrezId: 4292 EnsemblId: [ENSG00000076242](#)  
UniProt: [P40692](#) OMIM: [120436](#)



Variant 13: **Gene:** [PLG](#) **Your genotype:** C/T **Location:** chr6:161152905

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

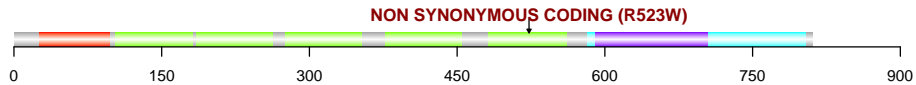
**Frequency:** **1KGenomes:** 0.00320 **dbSNP:** [rs4252129](#)

**Quality:** **Genotype quality:** 99 **Coverage depth:** 76

**Details:** **Gene description:** plasminogen  
**Transcript:** [ENST00000308192](#) **AA change:** R523W  
**EntrezId:** 5340 **EnsemblId:** [ENSG00000122194](#)  
**UniProt:** [P00747](#) **OMIM:** [173350](#)

PFAM (or SMART) domains for gene PLG, transcript ENST00000308192:

- PF00024: PAN-1\_domain
- PF00051: Kringle
- PF00089: Peptidase\_S1\_S6
- PF09342: Peptidase\_S1A\_nudel



Variant 14: **Gene:** [NEB](#) **Your genotype:** T/C **Location:** chr2:152435887

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

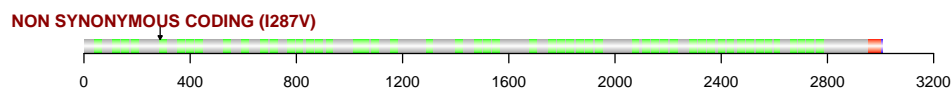
**Frequency:** **1KGenomes:** 0.00480 **dbSNP:** NA

**Quality:** **Genotype quality:** 99 **Coverage depth:** 54

**Details:** **Gene description:** nebulin  
**Transcript:** [ENST00000413693](#) **AA change:** I287V  
**EntrezId:** 4703 **EnsemblId:** [ENSG00000183091](#)  
**UniProt:** [P20929](#) **OMIM:** [161650](#)

PFAM (or SMART) domains for gene NEB, transcript ENST00000413693:

- PF00880: Nebulin\_35r-motif
- PF07653: SH3\_2
- PF00018: SH3\_domain

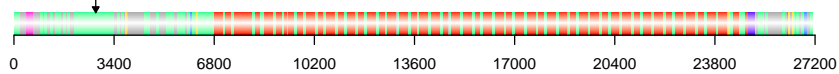


<b>Variant 15:</b>	<b>Gene:</b> <a href="#">TTN</a> <b>Your genotype:</b> C/A <b>Location:</b> chr2:179634961	<b>Type:</b> MODERATE
<b>Effect:</b>	<b>Impact:</b> NON SYNONYMOUS CODING	
<b>Frequency:</b>	<b>1KGenomes:</b> 0.00830	<b>dbSNP:</b> <a href="#">rs33917087</a>
<b>Quality:</b>	<b>Genotype quality:</b> 99	<b>Coverage depth:</b> 74
<b>Details:</b>	<b>Gene description:</b> titin <b>Transcript:</b> <a href="#">ENST00000342175</a> <b>EntrezId:</b> <a href="#">7273</a> <b>UniProt:</b> <a href="#">Q8WZ42</a>	<b>AA change:</b> V2777F <b>EnsemblId:</b> <a href="#">ENSG00000155657</a> <b>OMIM:</b> <a href="#">188840</a>

PFAM (or SMART) domains for gene TTN, transcript ENST00000342175:

- PF07679: Ig\_I-set
- PF09042: Titin\_Z
- PF00047: Immunoglobulin
- PF07686: Ig\_V-set
- PF00041: FN\_III
- PF00069: Se/Thr\_kinase-like\_dom
- PF07714: Ser-Thr/Tyr\_kinase

NON SYNONYMOUS CODING (V2777F)



# 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.