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:

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.
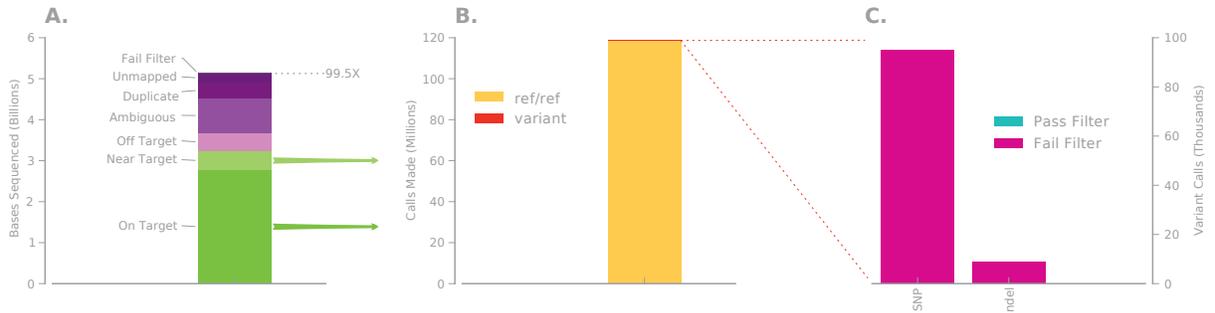
# Your exome in numbers



**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).
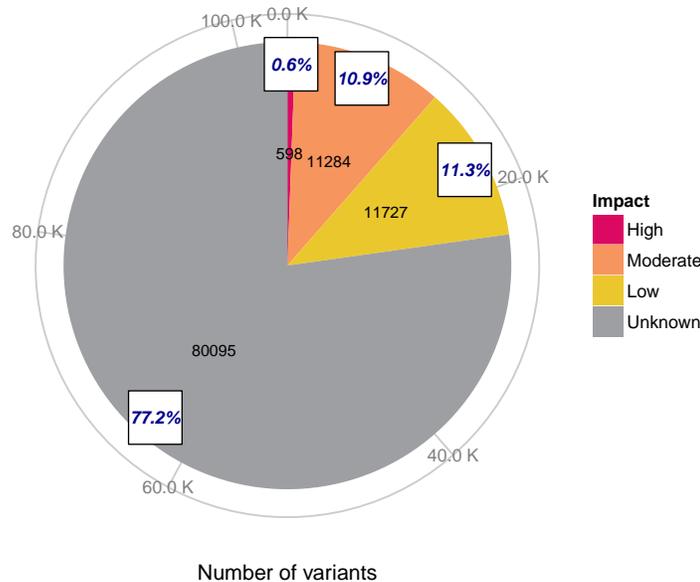
# Characterizing your variants



**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.
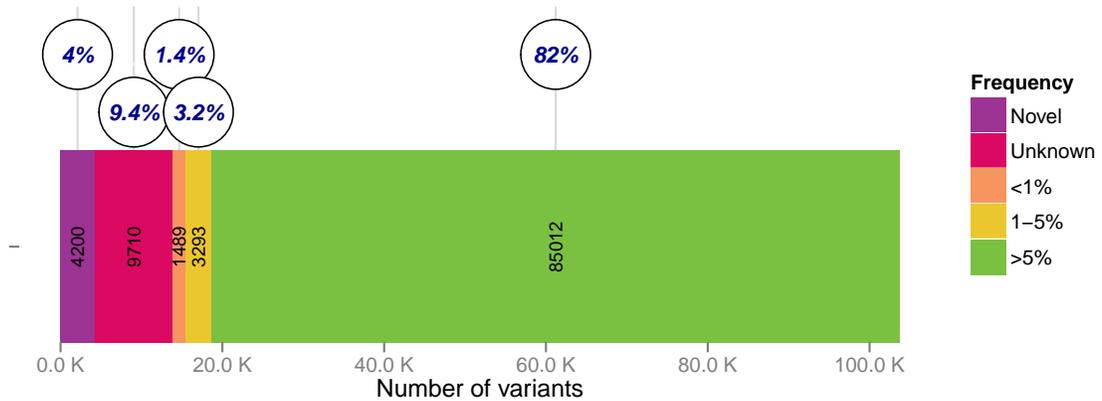
# How rare are your variants?



**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).
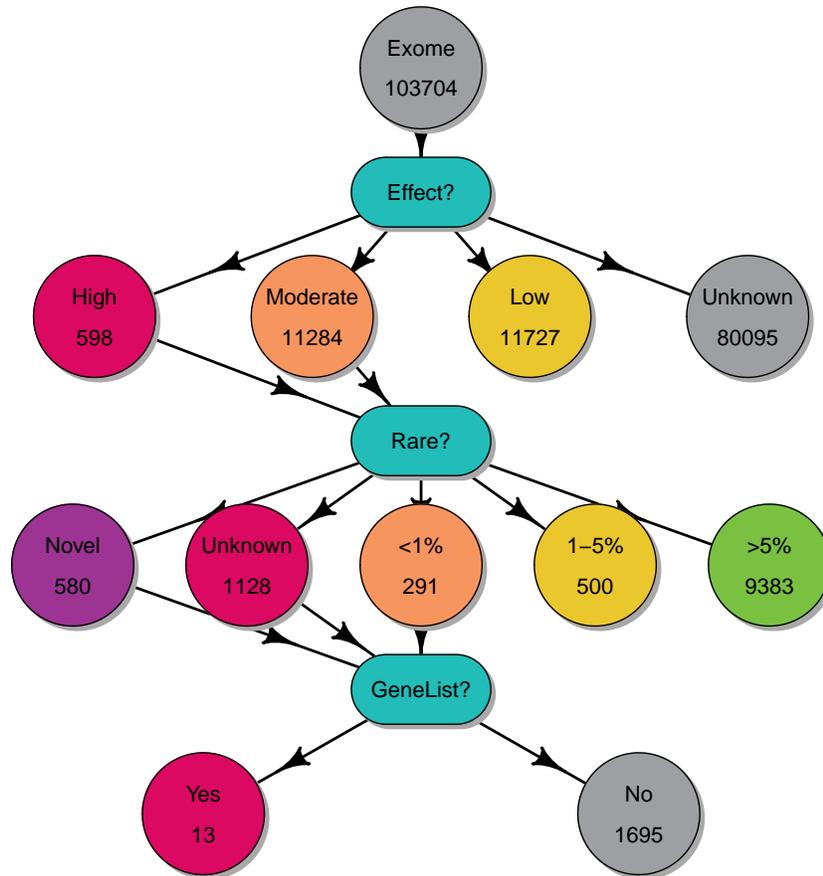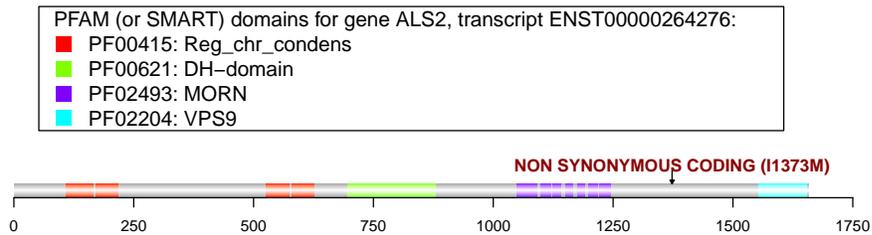
# Filtering your variants



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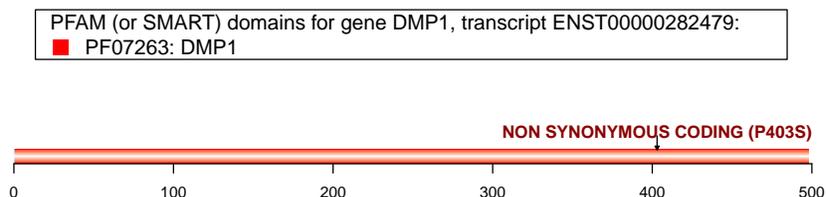
We hope you find this initial list of variants interesting and that it will help you in your journey through your exome. This short list of variants only scratches the surface of what your genome contains and is just the beginning of where your data can take you. Have fun!

# List of selected variants

| Variant 1: | Gene: ALS2 Your genotype: T/C Location: chr2:202575717 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00200 | dbSNP: rs61757691 |
| Quality: | Genotype quality: 99 | Coverage depth: 67 |
| Details: | Gene description: amyotrophic lateral sclerosis 2 (juvenile) | |
| | Transcript: ENST00000264276 | AA change: I1373M |
| | EntrezId: 57679 | EnsemblId: ENSG00000003393 |
| | UniProt: Q96Q42 | OMIM: 606352 |

PFAM (or SMART) domains for gene ALS2, transcript ENST00000264276:
- ■ PF00415: Reg_chr_condens
- ■ PF00621: DH–domain
- ■ PF02493: MORN
- ■ PF02204: VPS9

NON SYNONYMOUS CODING (I1373M)

0    250   500   750   1000   1250   1500   1750

| Variant 2: | Gene: DMP1 Your genotype: C/T Location: chr4:88584185 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 5e-04 | dbSNP: rs140275311 |
| Quality: | Genotype quality: 99 | Coverage depth: 84 |
| Details: | Gene description: dentin matrix acidic phosphoprotein 1 | |
| | Transcript: ENST00000282479 | AA change: P403S |
| | EntrezId: 1758 | EnsemblId: ENSG00000152592 |
| | UniProt: Q13316 | OMIM: 600980 |

PFAM (or SMART) domains for gene DMP1, transcript ENST00000282479:
- ■ PF07263: DMP1

NON SYNONYMOUS CODING (P403S)

0    100    200    300    400    500

| Variant 3: | Gene: ATXN2 Your genotype: T/C Location: chr12:111956226 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00200 | dbSNP: rs117851901 |
| Quality: | Genotype quality: 99 | Coverage depth: 63 |
| Details: | Gene description: ataxin 2 | |
| | Transcript: ENST00000535949 | AA change: N202S |
| | EntrezId: 6311 | EnsemblId: ENSG00000204842 |
| | UniProt: Q99700 | OMIM: 601517 |

PFAM (or SMART) domains for gene ATXN2, transcript ENST00000535949:
■ PF06741: LsmAD_domain
■ PF07145: Ataxin−2_C

**NON SYNONYMOUS CODING (N202S)**

| Variant 4: | Gene: CPT1A Your genotype: C/T Location: chr11:68562288 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 9e-04 | dbSNP: rs140958507 |
| Quality: | Genotype quality: 99 | Coverage depth: 52 |
| Details: | Gene description: carnitine palmitoyltransferase 1A (liver) | |
| | Transcript: ENST00000265641 | AA change: R288Q |
| | EntrezId: 1374 | EnsemblId: ENSG00000110090 |
| | UniProt: P50416 | OMIM: 600528 |

PFAM (or SMART) domains for gene CPT1A, transcript ENST00000265641:
■ PF00755: Carn_acyl_trans

**NON SYNONYMOUS CODING (R288Q)**

| Variant 5: | Gene: IGHMBP2 Your genotype: C/G Location: chr11:68702803 | |
|---|---|---|
| Effect: | Impact:   NON   SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 9e-04 | dbSNP: rs7122089 |
| Quality: | Genotype quality: 99 | Coverage depth: 28 |
| Details: | Gene description: immunoglobulin mu binding protein 2 | |
| | Transcript: ENST00000255078 | AA change: P557A |
| | EntrezId: 3508 | EnsemblId: ENSG00000132740 |
| | UniProt: P38935 | OMIM: 600502 |

PFAM (or SMART) domains for gene IGHMBP2, transcript ENST00000255078:
■ PF04851: UvrABC_suB
■ PF01424: R3H_ss–bd
■ PF01428: Znf_AN1

NON SYNONYMOUS CODING (P557A)

```
0     150    300    450    600    750    900    1050
```

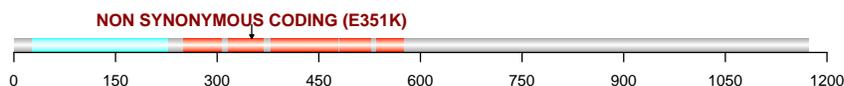| Variant 6: | Gene: LAMB3 Your genotype: C/T Location: chr1:209803163 | |
|---|---|---|
| Effect: | Impact:   NON   SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00180 | dbSNP: rs114875539 |
| Quality: | Genotype quality: 99 | Coverage depth: 64 |
| Details: | Gene description: laminin, beta 3 | |
| | Transcript: ENST00000356082 | AA change: E351K |
| | EntrezId: 3914 | EnsemblId: ENSG00000196878 |
| | UniProt: Q13751 | OMIM: 150310 |

PFAM (or SMART) domains for gene LAMB3, transcript ENST00000356082:
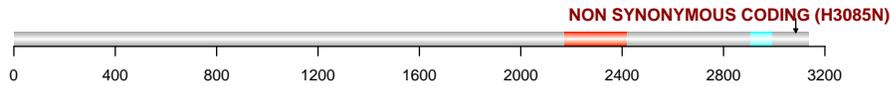■ PF00055: Laminin_N
■ PF00053: EGF_laminin

NON SYNONYMOUS CODING (E351K)

```
0     150    300    450    600    750    900    1050   1200
```

| **Variant 7:** | **Gene:** VPS13A **Your genotype:** C/A **Location:** chr9:80020874 | |
|---|---|---|
| **Effect:** | **Impact:** NON SYNONYMOUS CODING | **Type:** MODERATE |
| **Frequency:** | **1KGenomes:** 0.00640 | **dbSNP:** rs117983287 |
| **Quality:** | **Genotype quality:** 99 | **Coverage depth:** 69 |
| **Details:** | **Gene description:** vacuolar protein sorting 13 homolog A (S. cerevisiae) | |
| | **Transcript:** ENST00000376636 | **AA change:** H3085N |
| | **EntrezId:** 23230 | **EnsemblId:** ENSG00000197969 |
| | **UniProt:** Q96RL7 | **OMIM:** 605978 |

PFAM (or SMART) domains for gene VPS13A, transcript ENST00000376636:
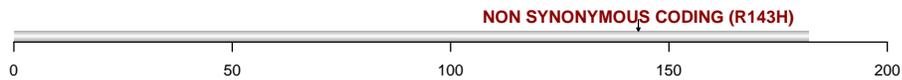- PF06650: VPSAP
- PF09333: Autophagy−rel_C

NON SYNONYMOUS CODING (H3085N)

0  400  800  1200  1600  2000  2400  2800  3200

| **Variant 8:** | **Gene:** MOGS **Your genotype:** C/T **Location:** chr2:74689078 | |
|---|---|---|
| **Effect:** | **Impact:** NON SYNONYMOUS CODING | **Type:** MODERATE |
| **Frequency:** | **1KGenomes:** 0.00470 | **dbSNP:** rs142032474 |
| **Quality:** | **Genotype quality:** 99 | **Coverage depth:** 40 |
| **Details:** | **Gene description:** mannosyl-oligosaccharide glucosidase | |
| | **Transcript:** ENST00000452063 | **AA change:** R507Q |
| | **EntrezId:** 7841 | **EnsemblId:** ENSG00000115275 |
| | **UniProt:** Q13724 | **OMIM:** 601336 |

PFAM (or SMART) domains for gene MOGS, transcript ENST00000452063:
- PF03200: Glycoside_hydrolase_63

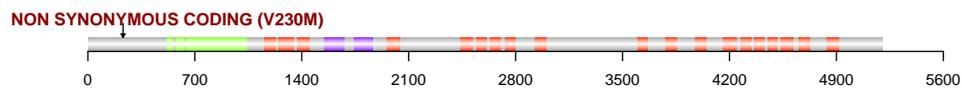NON SYNONYMOUS CODING (R507Q)

0  100  200  300  400  500  600  700  800

| Variant 9: | Gene: NPHP4 Your genotype: C/T Location: chr1:5951013 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00740 | dbSNP: rs34248917 |
| Quality: | Genotype quality: 99 | Coverage depth: 11 |
| Details: | Gene description: nephronophthisis 4 | |
| | Transcript: ENST00000378160 | AA change: R143H |
| | EntrezId: 261734 | EnsemblId: ENSG00000131697 |
| | UniProt: O75161 | OMIM: 607215 |

NON SYNONYMOUS CODING (R143H)

| 0 | 50 | 100 | 150 | 200 |
|---|---|---|---|---|

| Variant 10: | Gene: USH2A Your genotype: C/T Location: chr1:216538391 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00870 | dbSNP: rs45500891 |
| Quality: | Genotype quality: 99 | Coverage depth: 55 |
| Details: | Gene description: Usher syndrome 2A (autosomal recessive, mild) | |
| | Transcript: ENST00000307340 | AA change: V230M |
| | EntrezId: 7399 | EnsemblId: ENSG00000042781 |
| | UniProt: O75445 | OMIM: 608400 |

PFAM (or SMART) domains for gene USH2A, transcript ENST00000307340:
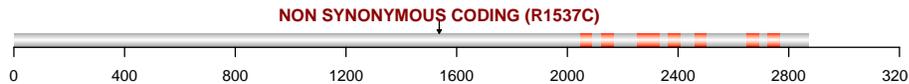- PF00053: EGF_laminin
- PF00041: FN_III
- PF00054: Laminin_G_1
- PF02210: Laminin_G_2

NON SYNONYMOUS CODING (V230M)

| 0 | 700 | 1400 | 2100 | 2800 | 3500 | 4200 | 4900 | 5600 |
|---|---|---|---|---|---|---|---|---|

| **Variant 11:** | **Gene:** DSP **Your genotype:** C/T **Location:** chr6:7581032 | |
|---|---|---|
| **Effect:** | **Impact:** NON SYNONYMOUS CODING | **Type:** MODERATE |
| **Frequency:** | **1KGenomes:** 0.00550 | **dbSNP:** rs28763967 |
| **Quality:** | **Genotype quality:** 99 | **Coverage depth:** 156 |
| **Details:** | **Gene description:** desmoplakin **Transcript:** ENST00000379802 **EntrezId:** 1832 **UniProt:** P15924 | **AA change:** R1537C **EnsemblId:** ENSG00000096696 **OMIM:** 125647 |

PFAM (or SMART) domains for gene DSP, transcript ENST00000379802:
■ PF00681: Plectin_repeat

NON SYNONYMOUS CODING (R1537C)

0    400    800    1200    1600    2000    2400    2800    320

| **Variant 12:** | **Gene:** NOTCH3 **Your genotype:** C/T **Location:** chr19:15273335 | |
|---|---|---|
| **Effect:** | **Impact:** NON SYNONYMOUS CODING | **Type:** MODERATE |
| **Frequency:** | **1KGenomes:** 0.00790 | **dbSNP:** rs115582213 |
| **Quality:** | **Genotype quality:** 99 | **Coverage depth:** 18 |
| **Details:** | **Gene description:** notch 3 **Transcript:** ENST00000263388 **EntrezId:** 4854 **UniProt:** Q9UM47 | **AA change:** V1952M **EnsemblId:** ENSG00000074181 **OMIM:** 600276 |

PFAM (or SMART) domains for gene NOTCH3, transcript ENST00000263388:
■ PF00008: EGF
■ PF07645: EGF_Ca−bd_2
■ PF07974: EGF_extracell
■ PF00066: Notch_dom
■ PF06816: Notch_NOD_dom
■ PF07684: Notch_NODP_dom
■ PF00023: Ankyrin_rpt
■ PF11936:

NON SYNONYMOUS CODING (V1952M)

0    300    600    900    1200    1500    1800    2100    2400

| Variant 13: | Gene: ALDH5A1 Your genotype: G/T Location: chr6:24505196 | |
|---|---|---|
| Effect: | Impact: NON SYNONYMOUS CODING | Type: MODERATE |
| Frequency: | 1KGenomes: 0.00820 | dbSNP: rs62621664 |
| Quality: | Genotype quality: 99 | Coverage depth: 32 |
| Details: | Gene description: aldehyde dehydrogenase 5 family, member A1 | |
| | Transcript: ENST00000546278 | AA change: A149S |
| | EntrezId: 7915 | EnsemblId: ENSG00000112294 |
| | UniProt: P51649 | OMIM: 610045 |

PFAM (or SMART) domains for gene ALDH5A1, transcript ENST00000546278:
■ PF00171: Aldehyde_DH_dom

NON SYNONYMOUS CODING (A149S)

0          100          200          300          400          500

# 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.