# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

PRAD: Prostate Adenocarcinoma

TCGA: The Cancer Genome Atlas

DNA: Deoxyribonucleic acid

RNA-Seq: RNA Sequencing

SNPs: Single Nucleotide Polymorphisms

CNV: Copy Nnmber Variation

Maf: Mutation Annotated Format

Vcf: Variant Calling Files

# Analysis and Annotation of exome sequencing data to identify and prioritize genes responsible for Prostate Adenocarcinoma

Ashish Chahal

Delhi Technological University

## ABSTRACT

After skin cancer prostate cancer is the second most prevalent cancer in men. Somatic mutations in Prostate Adenocarcinoma are revealed by processing of the next-generation DNA sequencing data of the exome region. Mutation in exome region directly effects the expression of the genes and sometimes inhibits the expression which can lead to several diseases. High throughput technologies and NGS analysis enable us to find out variations in the exome region that are involved in complex pathways of cancers. Biomarkers can be identified using NGS and exome sequencing analysis pipelines which can help in diagnosis, treatment and prognosis of the cancer. Exome play a major role in protein profiling so any change in this region affect the individual. PRAD exome data was used to analyze the variations in the exome region. Data for PRAD was downloaded from the TCGA web portal for tumor matched with normal types 17 samples on which exome sequence analysis pipeline were applied to predict and prioritize the genes involved for PRAD pathway.

Perl programming language was used to prioritize and analyze the exome data. Perl script maf2vcf.pl, DisGeNET, Annovar software packages were used to find out 93 probable genes that were filtered from DisGeNET. Then 54 genes were found in conserved regions with phastconselements46way score > 400.

17 TCGA IDs samples showed sequence alignment errors which were filtered by matching with segmented duplications. Polyphen2 annotations were used to give scores about the deleterious effect of the variants. After these steps we got the most probable genes that might be responsible for the cause of Prostate Adenocarcinoma (PRAD). GSTT1, TP53, CYP19A1, BRAF genes were already involved in the pathway of occurrence of prostate cancer and these genes were also present in the filtered genes in this study. Using experimental validation methods on the filtered genes we may help in finding out the novel genes that are involved in the complex pathway of prostate cancer.

# INTRODUCTION

Revolution in research is due to the wide range of applications and high throughput efficiency of the Next Generation Sequencing (NGS). NGS helps in studying the alternative splicing complexity (Martin *et al.,* 2011), landscape of mutations in cancer. The recent advancement in NGS following the first Human Genome draft (2003), massive data have been generated for various types of cancers. To analyze these data various computational and NGS pipelines approaches are used. The most widely used application is whole exome sequence analysis which is used to find the genetic basis of human diseases phenotype (Mamanova *et al.,* 2010). Many recent studies have been carrying out integrative analysis of epigenetic and exome sequences (Liao *et al.,* 2013).

Capture sequencing of the exome has been quite popular since such an approach showed its clinical utility to identify known and novel variants associated with Mendelian diseases. In a landmark report, exome capture followed by next generation sequencing was shown to have clinical relevance in Miller's syndrome. Furthermore, a number of researchers have used exome sequencing for diagnosis and identification of novel genetic variants and novels genes associated with Mendelian diseases. Additionally, it has also been used to finely map variations in complex diseases and used them in clinical diagnosis of cancers.

The Cancer Genome Atlas (TCGA) is a common platform designed to distribute and handle large volumes of research data for 34 types of cancer. Prostate Adenocarcinoma (PRAD) is a common type of cancer prevalent in North America, Australia, Europe and New Zealand and is a major cause of death in men (Hsing and Chokkalingam, 2006). TCGA provides information about the variations in exome data observed in the tumor samples in which somatic changes and variations are observed. The smart architecture of TCGA enables a researcher to download raw or processed data wherever applicable and available. Independent studies can be carried out to compare, analyze and interpret the information from various platforms on a particular sample. As is the mission of TCGA, the atlas of variations can be analyzed and stored to reduce the gap in the cancer and its molecular biology. Biomarkers can be found by using this study which helps in prognosis, diagnosis and treatment of the PRAD. These markers will be useful for personalized oncology treatment of the cancer.

Data analysis in this study, first involves the conversion of maf files into vcf files using perl language programming script and then several databases and packages were used to filter and analyze the data. In this study, the data was downloaded for PRAD from batch 184 for 17 tumor matched with normal samples.

In this study, we implemented ANNOVAR package (which uses the perl programming language), twobittofa tool to obtain the reference genome in correct format, DisGeNET database, maftovcf.pl perl scripts and MS-Excel were used to filter and analyze the variations in exome sequencing data. Finally the results were viewed on Integrative Genomics Viewer (IGV, Broad Institute).

# Review of literature

## 3.1 Next Generation Sequencing

The genome sequence is largely the same between individuals, all of us differ from each other in certain positions in the genome sequence. These changes are called genetic variations. These are associated with specific features which are shared with parents, so these variations are inherited one. The association of genetic variations and traits form the basis of human genetics. Now, not all genetic variations are associated with a human trait, but only a handful of them, mostly which fall in and around regions of protein coding genes (Conrad *et al.,* 2011). In many cases, the genetic variation is common in the population of individuals and these common variations are otherwise called polymorphisms.

Now sequencing individual genomes to understand the variations and arrive at the implications using the conventional Sanger sequencing methodology would have been extremely expensive. So the years which followed the human genome sequencing saw extensive investments into making nucleotide sequencing cheaper, fast and applicable in clinical settings. Consequent to these efforts, the field saw tremendous improvement in the throughput and speed and consequent drastic reduction in the costs of nucleotide sequencing. So much that it is now possible to sequence complete human genomes at a minuscule fraction of the cost incurred in the international human genome project. These technologies have been generally called next-generation sequencing technologies (NGS). As you would have imagined, NGS is not just one technology, but a generic name for a set of technologies which has enabled high-throughput sequencing of nucleotides. Each of the individual technologies significantly differ in the chemical reactions and readouts, but generally similar in the fact that they can sequence millions if not billions of sequences in one go (Koboldt *et al.,* 2010).

## 3.2 Exome Sequencing

The genome as we know is composed of over 3 billion bases. Not all of these three billion bases code for genes. Actually only a minuscule proportion of these 3 billion bases has protein coding potential. These regions are not contiguous in most cases of genes, but rather interspersed with regions which do not have a potential to encode for proteins. These are called exons. The exons encompass approximately over 50 million bases in the human genome. So one would naturally argue that it would be worthwhile to just sequence the 50 million odd bases in the human genome to understand genetic diseases. This is true in most cases with some exceptions. A number of mutations could also potentially affect the regulation and biogenesis of transcripts and could also cause diseases. Genetic mutations analysis show that these all fall in or near protein coding genes and exons.

Now sequencing the 50 million odd bases across the protein coding regions or exons and neighboring regulatory regions has been attempted. There are two popular approaches to 'capture' these regions. These are called capture methodologies. The principle is based on the fact that

complementary oligonucleotides would hybridize the region of interest and could be separated. Two popular capture methodologies have been developed; in first method we capture fragments on a glass surface, while the other captures it on magnetic beads. Fragments on capturing regions denature and then we do sequencing of these regions.
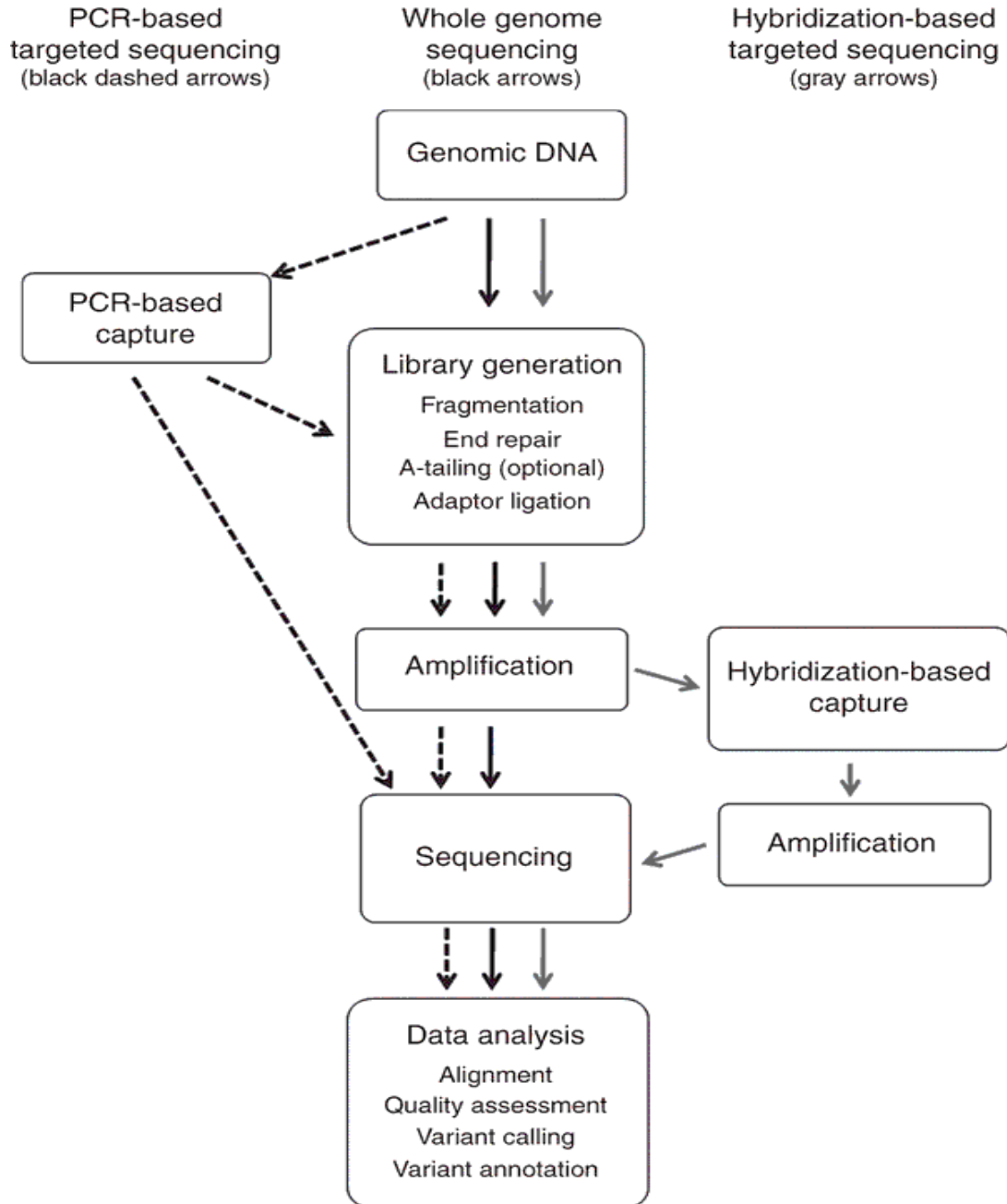


**Figure 1: NGS overview (Rehm *et al*., 2013)**

**3.3 Platform Design***:*

Data used in this study is shown in the table below that is downloaded from TCGA portal website (https://tcga-data.nci.nih.gov/tcga/tcgaPlatformDesign.jsp)

| Platform Center | TCGA Platform Code | | Platform Name |
|---|---|---|---|
| Broad.mit.edu | ILLUMINA GA_DNASeq | | Illumina Genome Analyzer DNA Sequencing |

**Table 2: Data center**

**3.4 The Cancer Genome Atlas (TCGA):**

The Cancer Genome Atlas stores 34 cancer types high level data that is generated by the various laboratories like Broad Institute, HGSC, WUSTL. These laboratories generate data from many samples of cancer patient in the form of clinical information, CNV(SNP arrays), DNA methylated data, micro RNA-seq data, somatic mutations, expression protein data, RNA seqV2 data, CNV(low Pass DNASeq), protected mutations data, copy number data. The data is present in different types of levels like Level 1, Level 2, Level 3 and each level has its own meaning when it comes to its data type. This data is present for tumor samples and also for normal samples. TCGA has unique ID for the samples.

In this study somatic mutation data is used that is downloaded from TCGA. Detail of data is below in the table.

TCGA provides data free available if it is not downloading able then you are not allowed to access that data because this data is protected. To access this data you need special permissions.

In TCGA data we have special bar codes for each sample. The specification of these bar codes is available at https://wiki.nci.nih.gov/display/TCGA/TCGA+barcode. Data can be downloaded for tumor and normal type samples. In this study I download **Level 2** data of prostate adenocarcinoma (PRAD) of **17 tumor matched with normal samples of batch 184**. Level 2 data is preprocessed data that is available on TCGA data portal. Level 3 data is not available as I mentioned it in the above table.

| Data Types | Cancer types applicable | Data type Name | Level 1 | Level 2 | Level 3 | Important Metadata | How to retrieve data |
|---|---|---|---|---|---|---|---|
| Mutations | All | Somatic mutation | Whole exome sequence | Somatic mutation for each participant | n/a | The mutation data do not have a standard MAGE-TAB archive associated with them yet. The latest mutated file specifications areavailable on the wiki. | Data Matrix &Bulk Download: Select 'Somatic Mutations for Data Type File Search: Select 'DNA Mutations for Data Category Publication MAF Search |

**Table 2*:* Data types and data levels at TCGA for Exome of various cancers**

9

**Figure 2: Data portal web face**

### 3.5 Analysis of exome level 2 data:

Exome Level 2 data of consists of the variations found in the exome region after preprocessing of the raw data and then match this data with reference genome. Then to find out or prioritizing the genes responsible for a particular disease here in this study DisGeNEt database and ANNOVAR package is used. After that annotation of these genes is also done using ANNOVAR package. Perl language is used to analyze level 2 data using ANNOVAR package and also perl scripts are used to covert maf files into vcf formats. EXCEL tools also used to filter genes from DisGeNET database.

ANNOVAR package is used to annotate functions of the variations up to date. Annovar provides various databases like refgene, cytoband, genomicSuperDups, Cosmic68, 1000g2012apr, esp6500, phastConslelements46way, snp138, ljb26.

Annovar package have perl scripts that uses these databases and helps in filtering those alterations that are responsible for the study of diseases. PhastContrastelements46way database helps in finding those alterations that fall in conserved regions. These variations are most likely responsible for disease. Then genomicSuperDups databases are used to filter out duplicated variations and it filtered our gene list further. Further databases are used to annotate these filtered genes.

### 3.6 Prostate adenocarcinoma

Prostate cancer is multi-factorial complex disease. Prostate cancer is more prevalent in North America, Australia, Europe and New Zealand (Jemal *et al.,* 2010). About 250,000 patients detected positively for prostate cancer in United States of America alone in 2015. However prostate cancer rate is slow in Asian countries. In India prostate cancer occurrence rate is very low 3.3/100,000(Shen *et al.,* 2010). In Indian cities Delhi, Pune, Kolkata and Thiruvananthapuram prostate cancer is the second leading cancer.  Mumbai and Bangaluru are the third leading sites. There is no standard clinical practice for prostate cancer yet however TP53( Balmukhanov *et al.,* 2013 ), GST family(Ecke *et al.,* 2010), CPY19A1(Kanda *et al*., 2015), PTEN and AR genes are probable genes responsible for prostate cancer.

### 3.7 Pipeline used in current study:

**(a)** Download TCGA data of exome in **maf** format

**(b)** Convert maf in to **vcf** format

**(c)** Then map these vcf files with the **DisGeNET** database

**(d)** Convert vcf files in to **avinput** format

**(e)** Filter genes by **phastConselements46way** database

**(f)** Further filtering of genes with **genomicSuperdups** database

**(g)** Validation of these filtered genes with **cosmics** database

**(h)** Gene annotation of these filtered genes

**(I)** In the end get the cytobands for the genes

# METHODOLOGY

## 4.1 Data Retrieval

Exome Data downloaded from TCGA data portal in the form of Somatic mutations that are in the form of BI Automated Mutation Calling for 17 same of tumor matched with normal provided with their TCGA IDs for Prostate Adenocarcinoma (PRAD).

| Data Type | Exome sequencing |
|---|---|
| Level | 2 |
| Center/Platform | Illumina Genome Analyzer DNA Sequencing |
| Batch | 184 |
| Disease | PRAD |

**Table 3: Data Detail for Prostate Adenocarcinoma (PRAD)**

Data downloaded in the maf format as look below



**Figure 3:  maf file**

## 4.2 Download Human Reference Genome hg19

Download the Human Reference Genome from the link (http://hgdownload.cse.ucsc.edu/downloads.html#human). Reference genome downloaded in the form of **bit** format.

To convert **bit** format into **fasta** format I use **twobittofa** utility tool of UCSC genome browser.

### 4.3 Converting maf files in to vcf files

Converting maf file into vcf files a perl script **maf2vcf.pl** is developed by ckandoth license to apache2. In this script a slight change is done manually by changing the location of the reference genome by giving the reference genome path that is in the user system. Then run the command.

### 4.4 Download the ANNOVAR software

ANNOVAR software is downloaded from the annovar web site but to download this software we have to register first and also needed an academic mail id. After registration a link is send to the mail id from where annovar can be easily downloaded.

### 4.5 Download DISGENET database

### 4.6 Downloaded the required databases that are needed in Annovar software

### 4.7 Convert the vcf files into .avinput format

Annovar uses only .avinput files so convert the vcf files into .avinput format by using convert2annovar.pl



**Figure 4 : convert2annovar.pl**

### 4.8 Gene based annotation of .avinput for extracting the names of the genes

To get the names of the genes that are filtered after conversion of maf files into vcf files is done by annotate_variation.pl

```
.exonic_variant_function
ashish@ashish-HP-Pavilion-15-Notebook-PC[annovar] clear                                                        [ 5:52PM]


ashish@ashish-HP-Pavilion-15-Notebook-PC[annovar] annotate_variation.pl -out 7211-01/7211 -build hg19 7211-01/7211.avinput humandb/   [ 5:54PM]
NOTICE: The --geneanno operation is set to ON by default
NOTICE: Reading gene annotation from humandb/hg19_refGene.txt ... Done with 50914 transcripts (including 11516 without coding sequence annotatio
n) for 26271 unique genes
NOTICE: Reading FASTA sequences from humandb/hg19_refGeneMrna.fa ... Done with 67 sequences
WARNING: A total of 345 sequences will be ignored due to lack of correct ORF annotation
NOTICE: Finished gene-based annotation on 45 genetic variants in 7211-01/7211.avinput
NOTICE: Output files were written to 7211-01/7211.variant_function, 7211-01/7211.exonic_variant_function
ashish@ashish-HP-Pavilion-15-Notebook-PC[annovar]                                                              [ 5:57PM]
ashish@ashish-HP-Pavilion-15-Notebook-PC[annovar] clear                                                        [ 6:03PM]
```

**Figure 6: .avinput script format**

Now run the same scripts for all 17 vcf files and then get the all genes name in the excel sheet.

Then mapping vcf files genes with the DisGeNET database.

This script annotate_variation.pl -out 7211-01/7211 -build hg19 7211-01/7211.avinput humandb/ output two files 7211-01/7211.variant_function and 7211-01/7211.exonic_variant_function.

After this step mapping .variant_function files genes name of the entire 17 sample we got filtered genes for PRAD.


## 4.9 Filtered variations in conserved region

To get the variation in conserved regions we have region based annotation to use the PhastConselements46way database which gives score and we filtered the variation above 400 threshold. This can be done by annotate_variation.pl -regionanno -build hg19 -out 7211-01/7211 -dbtype phastconselements46way 7211-01/7211.avinput humandb/.


## 4.10 Segmented duplicated variations filtered

This can be done first by downloading the genomicSuperDups database and then run the following perl script annotate_variation.pl -regionanno -build hg19 -out 7211-01/7211-dbype genomicSuperDups 7211-01/7211.avinput huamndb.

After this step I filtered out more genes from the above filtered gene (conserved region) by removing the variations from genomicSuperDups database.


## 4.11 Further filtration is done by polyphen2 annotation

Download HVAR database and then map the .avinput files with this database by annotate_variation.pl -filter -dbtype ljb23_pp2hvar -buildver hg19 -out 7211-01/7211 7211/7211.avinput humandb.

Probably damaging = .909 – 1

Possibly damaging = .447 - .908

Benign           = 0 - .446

**4.14 Mapping with cosmic database**

Cosmic database is Catalogue of Somatic mutation in cancers. Annovar uses this database and tell us that these mutations occur in cancer that is present in literature and how many times.

**4.15 Getting cytobands**

To get the location and the bands of variation download the cytoband database. Then uses the perl script annotate_variation.pl -regionanno -build hg19 -out 7211-01/7211 -dbtype cytoband 7211-01/7211.avinput humandb.

# RESULTS

## 5.1 Converting maf files into vcf files



**Figure 6: vcf file**

## 5.2 Converting vcf files into .avinput files
Annovar needs .avinput files for processing so we have to convert vcf into .avinput format.



**Figure 7: .avinput file**

## 5.3 DisGeNET database for Prostate Cancer



Figure 8: DisGeNET database

## 5.4 Filtering genes after mapping with DisGeNET database



**Figure 9: Matching vcf variants with DisGeNET**

## 5.5 Mapping with PhastConelements46way database

Filtering genes after matching with phastconselements46way we get the variations in conserved regions



| | Hugo_Sy | Entrez_Ge | Center | Ncbi_Buil | Chromoso | Start_Posi | End_Posit | Strand | Variant_C | Variant_T | Reference | Tumor_Se | Tumor_Se | Dbsnp_Rs | Dbsnp_Va | Tumor_Sa | Matched_ | Match_Nc | Match_Nc | Tumor_Va | Tumor_ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Hugo_Sy | Entrez_Ge | Center | Ncbi_Buil | Chromoso | Start_Posi | End_Posit | Strand | Variant_C | Variant_T | Reference | Tumor_Se | Tumor_Se | Dbsnp_Rs | Dbsnp_Va | Tumor_Sa | Matched_ | Match_Nc | Match_Nc | Tumor_Va | Tumor_ |
| 2 | AARS2 | 57505 | broad.mit. | 37 | chr6 | 44268380 | 44268380 | + | Silent | SNP | C | C | T | rs1426945 | by1000ger | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 3 | ABCA1 | 19 | broad.mit. | 37 | chr9 | 1.08E+08 | 1.08E+08 | + | Frame_Sh | DEL | AGAGGA | AGAGGA | - | | | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 4 | ABCA1 | 19 | broad.mit. | 37 | chr9 | 1.08E+08 | 1.08E+08 | + | Frame_Sh | DEL | AGAGGA | AGAGGA | - | | | TCGA-EJ | TCGA-EJ-7782-11A-01D-2114-08 | | | | |
| 5 | ABCA13 | 154664 | broad.mit. | 37 | chr7 | 48352729 | 48352729 | + | Silent | SNP | C | C | T | | | TCGA-H( | TCGA-HC-7742-11A-01D-2114-08 | | | | |
| 6 | ABCA13 | 154664 | broad.mit. | 37 | chr7 | 48352729 | 48352729 | + | Silent | SNP | C | C | T | | | TCGA-H( | TCGA-HC-7742-10A-01D-2115-08 | | | | |
| 7 | ABCC11 | 85320 | broad.mit. | 37 | chr16 | 48212570 | 48212570 | + | Missense_ | SNP | G | G | A | | | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 8 | ABCC5 | 10057 | broad.mit. | 37 | chr3 | 1.84E+08 | 1.84E+08 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7331-11A-01D-2114-08 | | | | |
| 9 | ABCC5 | 10057 | broad.mit. | 37 | chr3 | 1.84E+08 | 1.84E+08 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7331-10A-01D-2114-08 | | | | |
| 10 | ABCC5 | 10057 | broad.mit. | 37 | chr3 | 1.84E+08 | 1.84E+08 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 11 | ABCC5 | 10057 | broad.mit. | 37 | chr3 | 1.84E+08 | 1.84E+08 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-11A-01D-2114-08 | | | | |
| 12 | ABCC8 | 6833 | broad.mit. | 37 | chr11 | 17428224 | 17428224 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 13 | ABCC8 | 6833 | broad.mit. | 37 | chr11 | 17428224 | 17428224 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-11A-01D-2114-08 | | | | |
| 14 | ABCE1 | 6059 | broad.mit. | 37 | chr4 | 1.46E+08 | 1.46E+08 | + | Missense_ | SNP | A | A | C | | | TCGA-H( | TCGA-HC-7745-10A-01D-2115-08 | | | | |
| 15 | ABCE1 | 6059 | broad.mit. | 37 | chr4 | 1.46E+08 | 1.46E+08 | + | Missense_ | SNP | A | A | C | | | TCGA-H( | TCGA-HC-7745-11A-01D-2114-08 | | | | |
| 16 | AC002331 | 0 | broad.mit. | 37 | chr16 | 26599065 | 26599068 | + | RNA | DEL | ACAG | ACAG | - | rs71134607 | | TCGA-H( | TCGA-HC-7752-10A-01D-2115-08 | | | | |
| 17 | AC006050 | 0 | broad.mit. | 37 | chr17 | 28901676 | 28901676 | + | RNA | DEL | C | C | - | | | TCGA-EJ | TCGA-EJ-7317-10A-01D-2114-08 | | | | |
| 18 | AC006050 | 0 | broad.mit. | 37 | chr17 | 28901676 | 28901676 | + | RNA | DEL | C | C | - | | | TCGA-EJ | TCGA-EJ-7317-11A-01D-2114-08 | | | | |
| 19 | AC008103 | 0 | broad.mit. | 37 | chr22 | 18844766 | 18844766 | + | RNA | SNP | G | G | A | | | TCGA-EJ | TCGA-EJ-7328-10A-01D-2114-08 | | | | |
| 20 | AC010547 | 0 | broad.mit. | 37 | chr16 | 71516014 | 71516014 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-10A-01D-2114-08 | | | | |
| 21 | AC010547 | 0 | broad.mit. | 37 | chr16 | 71516014 | 71516014 | + | Missense_ | SNP | C | C | T | | | TCGA-EJ | TCGA-EJ-7782-11A-01D-2114-08 | | | | |
| 22 | AC015818 | 0 | broad.mit. | 37 | chr17 | 20424245 | 20424245 | + | RNA | SNP | C | C | T | | | TCGA-H( | TCGA-HC-7819-11A-01D-2114-08 | | | | |
| 23 | AC018730 | 0 | broad.mit. | 37 | chr2 | 1.05E+08 | 1.05E+08 | + | RNA | DEL | TGGTGA | TGGTGA | - | rs1509372 | by1000ger | TCGA-EJ | TCGA-EJ-7794-11A-01D-2114-08 | | | | |
| 24 | AC019118 | 0 | broad.mit. | 37 | chr2 | 2910378 | 2910379 | + | RNA | INS | - | - | C | rs7205115 | by1000ger | TCGA-EJ | TCGA-EJ-7317-11A-01D-2114-08 | | | | |
| 25 | AC019118 | 0 | broad.mit. | 37 | chr2 | 2910768 | 2910769 | + | RNA | INS | - | - | T | rs1425859 | by1000ger | TCGA-H( | TCGA-HC-7737-11A-02D-2114-08 | | | | |

**Figure 10: Conserved region variations**

**5.6 list of genes after mapping with Phastconselements Database**

Variations in the genes are filtered on the basis of the threshold score above 400. Variations above this score are likely to fall in conserved region.

ANKHD1, APC, ASAH1, ASCC2, ASXL1, ATM, BRAF, BUD13, CLTC, CXCR2, **CYP19A1**, DNAH8, EIF4G2, FGFR2, FLT4, **FOXA1**, FOXP2, FRS3, FSTL1, FZR1, GPX1, **GSTT1**, HDAC9, HPCA, KDR, KRIT1, LPAR1, MED15, MMP14, MOAP1, NCOA3, NCOR2, NOTCH1, PLK3, PLK4, PMEPA1, PPIG, RAF1, RB1, RCHY1, REPS1, SCN9A, SEMA3A, SIRT1, SLC39A6, STK26, TCN1, TGFBI, THBS1, **TP53**, TRAM1, TSHZ3, TTN.

**5.7 Region based annotation with genomicSuperDups**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| genomicSuperDups | Score=0.944882;Name=chr13:19239253 | chr19 | 4791992 | 4791992 | A | - | het |
| genomicSuperDups | Score=0.967305;Name=chr9:68714810 | chr20 | 3E+07 | 3E+07 | T | C | het |
| genomicSuperDups | Score=0.950851;Name=chr21:46011302 | chr21 | 4.6E+07 | 4.6E+07 | CTGCTG | - | het |
| genomicSuperDups | Score=0.964372;Name=chr1:144614531 | chr1 | 1.7E+07 | 1.7E+07 | C | T | het |
| genomicSuperDups | Score=0.98083;Name=chr1:148271039 | chr1 | 1.5E+08 | 1.5E+08 | T | A | het |
| genomicSuperDups | Score=0.98083;Name=chr1:148271039 | chr1 | 1.5E+08 | 1.5E+08 | A | G | het |
| genomicSuperDups | Score=0.948211;Name=chr5:140553661 | chr5 | 1.4E+08 | 1.4E+08 | AGGCCG | - | het |
| genomicSuperDups | Score=0.969561;Name=chr13:25153561 | chr13 | 2.6E+07 | 2.6E+07 | - | AAAAAC | het |
| genomicSuperDups | Score=0.902673;Name=chr7:142121806 | chr7 | 1.4E+08 | 1.4E+08 | AGG | - | het |
| genomicSuperDups | Score=0.995648;Name=chr15:32445406 | chr15 | 3E+07 | 3E+07 | - | T | het |

**Figure 11: GenomicSuperDups filtration**

After filtering with genomicSuperDups segmented duplicated genetic variants can be removed. These types of variants can be show two non polymorphic types in genome.
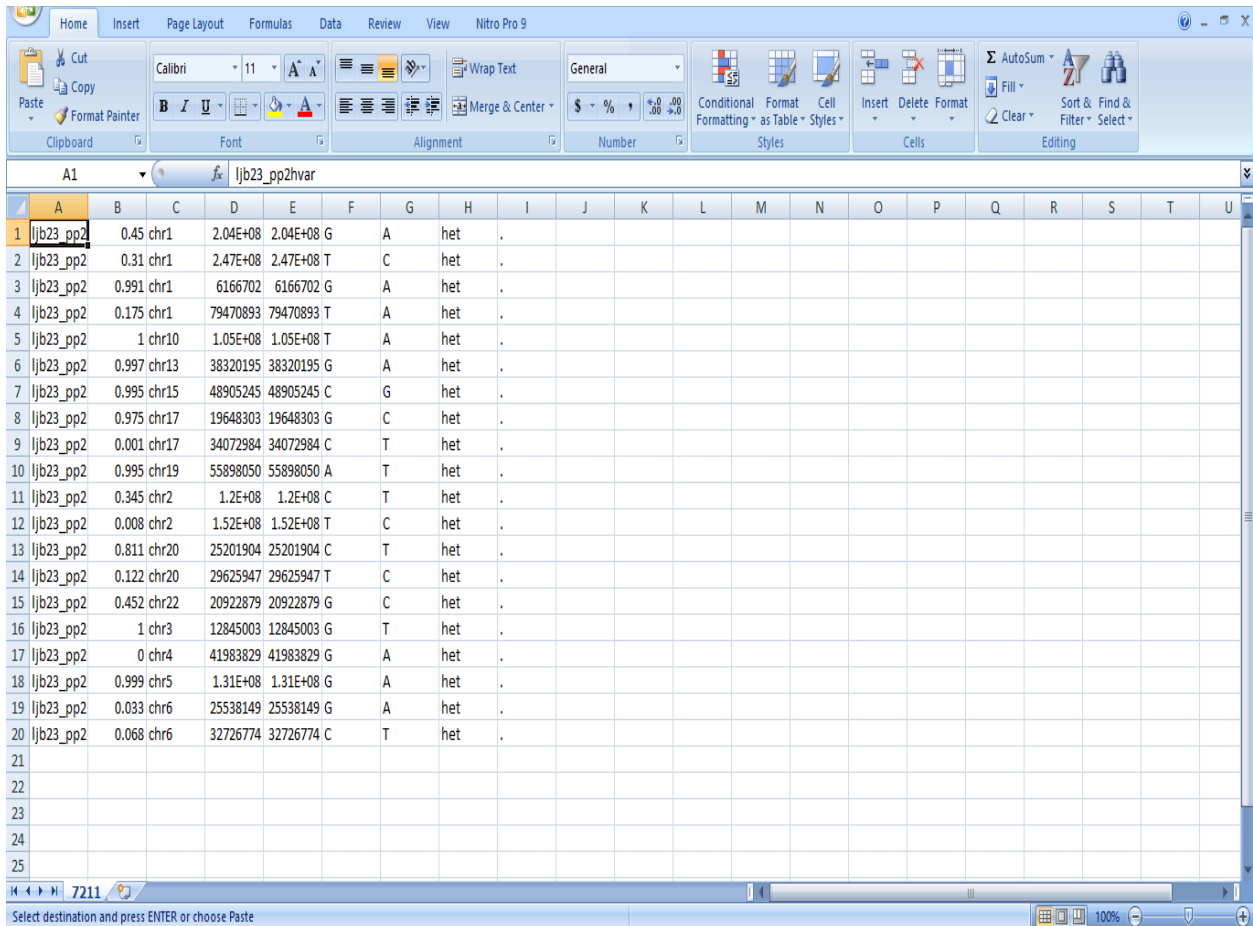
## 5.8 Filter based annotation with polyphen2 database to know about the damaging effect of disease

This annotation shows the damaging effect of the variants on the basis of the scores.

Probably damaging = .909 – 1

Possibly damaging = .447 - .908

Benign = 0 - .446



| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | ljb23_pp2 | 0.45 | chr1 | 2.04E+08 | 2.04E+08 | G | A | het | . |
| 2 | ljb23_pp2 | 0.31 | chr1 | 2.47E+08 | 2.47E+08 | T | C | het | . |
| 3 | ljb23_pp2 | 0.991 | chr1 | 6166702 | 6166702 | G | A | het | . |
| 4 | ljb23_pp2 | 0.175 | chr1 | 79470893 | 79470893 | T | A | het | . |
| 5 | ljb23_pp2 | 1 | chr10 | 1.05E+08 | 1.05E+08 | T | A | het | . |
| 6 | ljb23_pp2 | 0.997 | chr13 | 38320195 | 38320195 | G | A | het | . |
| 7 | ljb23_pp2 | 0.995 | chr15 | 48905245 | 48905245 | C | G | het | . |
| 8 | ljb23_pp2 | 0.975 | chr17 | 19648303 | 19648303 | G | C | het | . |
| 9 | ljb23_pp2 | 0.001 | chr17 | 34072984 | 34072984 | C | T | het | . |
| 10 | ljb23_pp2 | 0.995 | chr19 | 55898050 | 55898050 | A | T | het | . |
| 11 | ljb23_pp2 | 0.345 | chr2 | 1.2E+08 | 1.2E+08 | C | T | het | . |
| 12 | ljb23_pp2 | 0.008 | chr2 | 1.52E+08 | 1.52E+08 | T | C | het | . |
| 13 | ljb23_pp2 | 0.811 | chr20 | 25201904 | 25201904 | C | T | het | . |
| 14 | ljb23_pp2 | 0.122 | chr20 | 29625947 | 29625947 | T | C | het | . |
| 15 | ljb23_pp2 | 0.452 | chr22 | 20922879 | 20922879 | G | C | het | . |
| 16 | ljb23_pp2 | 1 | chr3 | 12845003 | 12845003 | G | T | het | . |
| 17 | ljb23_pp2 | 0 | chr4 | 41983829 | 41983829 | G | A | het | . |
| 18 | ljb23_pp2 | 0.999 | chr5 | 1.31E+08 | 1.31E+08 | G | A | het | . |
| 19 | ljb23_pp2 | 0.033 | chr6 | 25538149 | 25538149 | G | A | het | . |
| 20 | ljb23_pp2 | 0.068 | chr6 | 32726774 | 32726774 | C | T | het | . |

Figure 12: Polyphen2 scores

## 5.9 Mapping input file with cosmic database

Validation of the filtered genes can be done by this database. Cosmic database contains all the somatic mutations that present in literature of the all type of the cancers.



**Figure 13: Cosmic validation**

## 5.10 Gene annotation to know about the nature of the variant and its effect on amino acid change



**Figure 14: Exonic variants**

## 5.11 Predict the variations lie in the Transcription factor binding sites
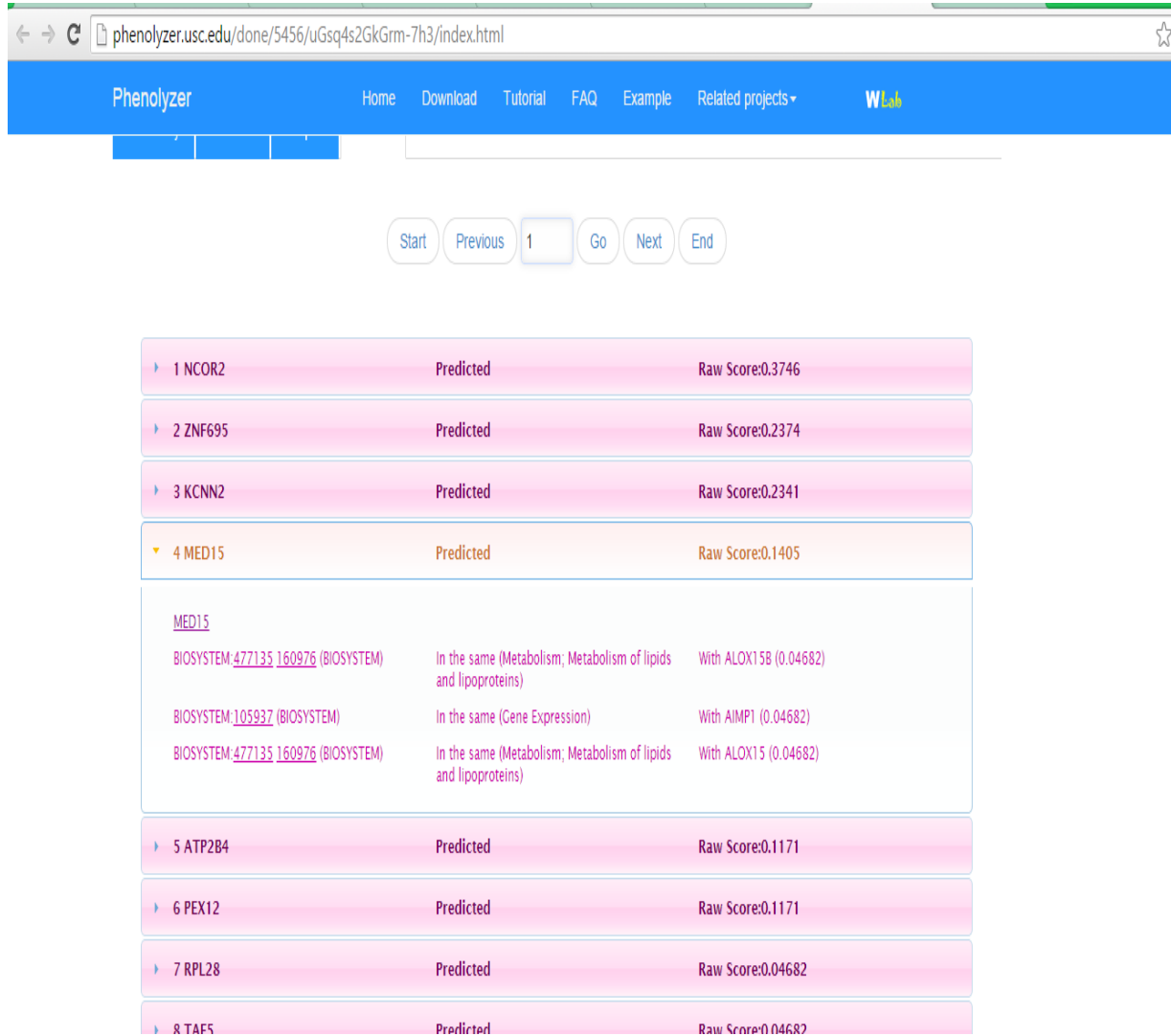


**Figure 15: Transcription factor binding site**

## 5.12 Getting cytobands

Chromosomes bands that are Geimsa- stained can be achieved by mapping with the cytobands database



**Figure 16: CytoBand**

**5.13 Phenolyzer result:** Phenolyzer analysis results show the network and the pathways of the genes in which they involve.



**Figure 17: Phenolyzer Details**

# DISCUSSION

With the advancement in NGS, there is sudden explosion of data in the public centers (Rakyan *et al.,* 2011). Genetic alterations accumulating in the exome region of the genome are the main reason for the occurrence of various forms of cancers (Meyerson *et al*., 2010). Since variations in the exome region directly affect the phenotype of the individuals (Pauline *et al*., 2008), many human diseases involve genetic alterations in the exome region such as cancer, Mendelian disorders and neurological diseases.

Exome sequence analysis of the genetic variations in PRAD patients were carried out in this study. A number of genes were filtered after analysis that carries those variations that might be the cause of the PRAD. From these filtered genes few genes like GSST1, TP53 and CYP19A1 are already known to play the role in PRAD pathway. Similarly, we could conclude that out of the 52 filtered genes some could possibly act as biomarkers.

Segmented duplicated variations that show two non-polymorphic phenotypes were removed from the filtered genes and this was done by matching with genomicSuperDups database. To know about the damaging effect of the variation, polyphen2 was used (Adzhubei *et al.,* 2010). The variations with scores in the range 0.446-0 are benign, score ranges 0.447-0.908 are possibly damaging and score 0.909-1.0 are probably damaging. 52 Filtered genes variations had amino acid changes which were annotated using Refgene database. Type of mutations like frame shift, synonymous and non synonymous were also found out using the same database. Database tfbsConsSites showed that which variants lie in the conserved regions of the transcription factor binding sites on the basis of scores. Cytogenetic bands of each gene were obtained by matching with cytoband database and viewed using Integrative Genomics Viewer. Network pathways of each filtered gene were also annotated. This analysis tells about the Phenolyzer effects of the genes like genes involved in lipid metabolism etc. Validation of the filtered genes was done by COSMIC database. This database includes somatic mutations reported in the literature. This study can be used on cancer and other type of diseases to find out the probable biomarkers.

To validate the results, wet lab experiments are very important for supporting the analysis of genetic variations in the exome regions. After validation of the results, variations can be used as biomarkers for the prostate cancer which could help in prognosis and diagnosis of the cancer. Different phenotypic effects are observed in the same gene that shows different mutations (Boyadjiev *et al.,* 2000). Environmental factors also affect some phenotypes due to the presence of some variations (Hunter, 2005). So along with genetic factors we should also consider the environmental factors for the causation of the cancer.

# CONCLUSION

Somatic mutations in Prostate Adenocarcinoma are revealed by processing of the next-generation DNA sequencing data of the exome region. NGS helps in studying the alternative splicing complexity (Martin *et al*., 2011), landscape of mutations in cancer. Exome sequencing used for diagnosis and identification of novel genetic variants and novels genes associated with Mendelian diseases.

In present study 93 genes filtered by mapping vcf files genetic variations with the DisGeNEt database. Further filtration was done by matching with the phastconselelement46way and genomicSuperDups database. After this step we got 52 most probable variations that might alter the pathways that caused prostate cancer.

Annotation of these variations was done by ANNOVAR software package. Annotation includes to find the nature of the variants; like: are variants falling in transcription binding factor sites, class of mutation in which variants lie, to know about the damaging effects of the variants and to get the cytobands of the variations. These annotations use the tnbps, refgene, polyphen2 and cytoband databases respectively. Validation of the filtered genes can be done by matching the variations with the Cosmic database, this database consists of the variations that present in the literature. In this study annotations include gene based, region based and filter base analysis.

Advancement in integrative approaches like somatic mutations analysis with DNA methylations and gene expression studies, we could expect a well developed personalized medicine in future soon (Rabbani *et al.,* 2014).

# FUTURE PROSPECTIVE

In this study a computational pipeline is developed for exome sequence analysis to prioritize and annotate the genetic variations found in exome region. Hg19 Human Reference Genome was used for the initial analysis work in this study. This study marks as a proof of concept for integrating the studies of different platforms such as DNA methylation, somatic mutations, gene expressions, and clinical data along with exome sequencing. TCGA has enabled the researchers to carry out a multi-platform analysis as it provides multi-platform data for the samples of same patient. TCGA is maintaining data very well and the architecture is well documented.

In the current study analysis and annotation of somatic mutations obtained by exome sequencing was done by ANNOVAR software. This is another progress for the further development of the pipelines for the analysis of the exome data. Integrative studies are also done on somatic variations with the DNA methylated and gene expression studies. Advancement in NGS and collaboration with statisticians, mathematicians, computer engineers and bioinformaticians could help in better development of the new pipelines that can integrate somatic variations, clinical data, DNA expression, DNA methylated and SNP data to find out the functional biomarkers. These biomarkers can be used for the diagnosis, prognosis and treatment of the complex cancers.

Whole exome sequencing can be used for analysis of other diseases like monogenic disorders, hearing loss, intellectual disabilities, movement disorders, cardiovascular diseases, obesity and diabetes, cancers and hypertensions.

More accurate Bioinformatics tools will be required for studying genetic diseases so that we can counsel the patient, find the more precise diagnosis and treatment for the disease in a personalized manner. This study will also further contribute to the development of personalized medicine. In future, as the consequence of advancement in integrative approaches like somatic mutations analysis with DNA methylations and gene expression studies, we could expect a well developed personalized medicine in future soon.

There are ethical issues also concerned in medical genetic studies. While studying the individual genome, it can reveal the relative DNA information or even exposed the properties of the population. Confidentiality, privacy and return of the results are some ethical issues that should be solved by professionals.

# REFERENCES

Martin, JA; Wang, Z (2011). Next-generation transcriptome assembly. Nat Rev Genet. 12 (10):671–682

Vogelstein, B; Papadopoulos, N; Velculescu, VE; Zhou, S; Diaz, LA Jr; Kinzler, KW (2013). Cancer genome landscapes. Science. 339 (6127):15461558

Conrad, DF; Keebler, JE; DePristo, MA; Lindsay, SJ; Zhang, Y; Casals, F; Idaghdour, Y; Hartl, CL; Torroja, C; Garimella, KV; Zilversmit, M; Cartwright, R; Rouleau, GA; Daly, M; Stone, EA; Hurles, ME; Awadalla, P; 1000 Genomes Project (2011). Variation in genome-wide mutation rates within and between human families. Nat Genet. 43(7):712-4

Wang, K; Li, M; Hakonarson, H (2010). ANNOVAR; functional annotation of genetic variants from high throughput sequencing data. Nucleic Acid Res. 38(16).

Perkel, MJ (2010). BIG TROUBLE IN LITTLE SCIENCE Nature Methods. 7: 589–593

Jemal, A: Bray, F; Center, MM; Ferlay, J; War, E; Forman, D (2011). Global cancer statistics. CA Cancer J Clin. 61(2):69-90

Shen, MM; Abate, C (2010). Molecular genetics of prostate cancer: new prospects for old challenges.Genes Dev. GD 24(18):1967-2000

Balmukhanov, T; Khanseitova, A; Nigmatova, V; Ashirbekov, E; Talaeva, S; Aitkhozhina, N. (2013) Polymorphisms at *GSTM* 1, *GSTP* 1, *GSTT* 1 Detoxification Genes Loci and Risk of Breast Cancer in Kazakhstan Population. Advances in Breast Cancer Research. 2(4):114-118

Mamanova, L; Coffey, AJ; Scott, CE; Kozarewa, I; Turner, EH; Kumar, A; Howard, E; Shendure, J; Turner, DJ (2010). Target-enrichment strategies for next-generation sequencing. Nat Methods. 7(2):111-118

Damodaran, S; Berger, MF; Roychowdhury, S (2015). Clinical tumor sequencing: opportunities and challenges for precision cancer medicine. Europe PubMed Central. 35: e175-182

M, Meyerson; S, Gabriel; G, Getz(2010). Advances in understanding cancer genomes through second-generation sequencing. Nature Reviews Genetics. 11:685-696

Pauline,C. Ng ; Samuel, Levy; Jiaqi, Huang; Timothy, B; Stockwell; Brian, P. Walenz; Kelvin, Li; Nelson, Axelrod; Dana, A; Busam, Robert, L; Strausberg, J (2008). Genetic Variation in an Individual Human Exome. PLoS Genet. 4(8): e1000160

Adzhubei, IA; Schmidt, S; Peshkin, L; Ramensky, VE; Gerasimova, A; Bork, P; Kondrashov, AS; Sunyaev, SR (2010). Annotation of functional variation in personal genomes using RegulomeDB. Nat Methods. 7(4):248-249

Boyadjiev, SA; Jabs, EW(2000). Online Mendelian Inheritance in Man (OMIM) as a knowledge base for human developmental disorders. Clin Genet. 57(4):253-66

Hunter, DJ(2005). Gene-environment interactions in human diseases. Nat Rev Genet. 6(4):287-98.
Rakyan VK, Down TA, Balding DJ, *et al.* (2011). Epigenome-wide association studies for common human diseases. Nat Rev Genet; 12:529–41.

Spinelli, A; R, Piazza; Pirola, A; Valletta, S (2014). Whole-Exome Sequencing Data–Identifying Somatic Mutations. Springer Link. 419-427

Rabbani, B; Mahdieh, N; Hosomichi, K; Nakaoka, H; Inoue, I(2012). Next-generation sequencing: impact of exome sequencing in characterizing Mendelian disorders. J. Hum. Genet. **57**, 621–632

B, Rabbani; M, Tekin; N, Mahdieh (2014). The promise of whole-exome sequencing in medical genetics. Journal of Human Genetics. 59: 5–15

Biesecker, L; Mullikin, J. C; Facio, F. M; Turner, C; Cherukuri, P. F; Blakesley, R. W(2009). The ClinSeq Project: piloting large-scale genome sequencing for research in genomic medicine. Genome Res. 19: 1665–1674

Ecke, TH1; Schlechte, HH; Schiemenz, K; Sachs, MD, Lenk, SV; Rudolph, BD; Loening, SA(2010). TP53 gene mutations in prostate cancer progression. Anticancer Res. 30(5):1579-86

Kanda, S; Tsuchiya, N; Narita, S; Inoue, T; Huang, M; Chiba, S; Akihama, S; Saito, M; Numakura, K; Tsuruta, H; Satoh, S; Saito, S; Ohyama, C; Arai, Y; Habuchi, T(2015). Effects of functional genetic polymorphisms in the CYP19A1 gene on prostate cancer risk and survival. Int J Cancer. 136(1):74-82

# APPENDIX

1. annotate_variation.pl -buildver hg19 -downdb -webfrom annovar refGene humandb/
2. annotate_variation.pl -buildver hg19 -downdb cytoBand humandb/
3. annotate_variation.pl -buildver hg19 -downdb genomicSuperDups humandb/
4. annotate_variation.pl -buildver hg19 -downdb -webfrom annovar esp6500siv2_all humandb/
5. annotate_variation.pl -buildver hg19 -downdb -webfrom annovar 1000g2014oct humandb/
6. annotate_variation.pl -buildver hg19 -downdb -webfrom annovar snp138 humandb/
7. annotate_variation.pl -buildver hg19 -downdb -webfrom annovar ljb26_all humandb/
8. convert2annovar.**pl** -**format** vcf4 example/ex2.vcf > ex2.avinput
9. annotate_variation.pl -**out** ex1 -build hg19 example/ex1.avinput humandb/
10. annotate_variation.pl -build hg19 -downdb phastConsElements46way humandb/
11. annotate_variation.pl -regionanno -build hg19 -out ex1 -dbtype phastConsElements46way example/ex1.avinput humandb/
12. annotate_variation.**pl** -build hg19 -downdb tfbsConsSites humandb/
13. annotate_variation.**pl** -regionanno -build hg19 -**out** ex1 -dbtype tfbsConsSites example/ex1.avinput humandb/
14. annotate_variation.pl -build hg19 -downdb cytoBand humandb/
15. annotate_variation.pl -regionanno -build hg19 -**out** ex1 -dbtype cytoBand example/ex1.avinput humandb/
16. annotate_variation]$ annotate_variation.pl -build hg19 -downdb genomicSuperDups humandb/
17. annotate_variation.pl -regionanno -build hg19 -out ex1 -dbtype genomicSuperDups example/ex1.avinput humandb/
18. annotate_variation.pl -filter -dbtype ljb23_pp2hvar -buildver hg19 -**out** ex1 example/ex1.avinput humandb/