Annotating and prioritizing SNVs

Practical

Pre-requisites

- Check your folder for annovar, we will need the following scripts:
 - convert2annovar.pl
 - annotate_variation.pl
 - table_annovar.pl
- There should be a database folder humandb that should contain the following databases (Note: these have been downloaded for you with annovar_commands.sh):
- You should have your filtered vcf file from the mutation caller ready

Preparation

- Annovar uses their own format for the input files
- Generate this file from the vcf using:

```
$ software/annovar/convert2annovar.pl \
    -format vcf4old \
    path_to_your_input.vcf.gz \
    > path_to_your_output.avinput;
```

Output:

```
C02Q11SYFVH6:CRUK_summerschool perner01$ head tmp/HCC1143_vs_HCC1143_BL.annot.muts.avinput
        10150
                10150
                                          unknown .
                                                           449
        10180
                10180
                                                           270
                                          unknown .
        10241
                10241
                                          unknown .
                                                           172
        10291
                10291
                         C
                                          unknown .
                                                           191
                10315
                         C
        10315
                                                           251
                                          unknown .
                10348
        10348
                                                           395
                                          unknown .
        10354
                10354
                                                           574
                                          unknown .
        10357
                10357
                                          unknown .
                                                           560
        10394
                10394
                                                           624
                                          unknown .
        10440
                                                           600
                10440
                                          unknown .
```

Gene-based annotation

- Annovar performs gene-based annotation as default
- Will generate at once annotation
 - with respect to genes and
 - with respect to functional effect on coding sequence

```
$ software/annovar/annotate_variation.pl \
    --buildver hg19 \
    path_to_your_annovar_file.avinput \
    path_to_your_db_folder;
```

Output:

```
C02Q11SYFVH6:CRUK_summerschool perner01$ ls tmp/
HCC1143_vs_HCC1143_BL.annot.muts.avinput
HCC1143_vs_HCC1143_BL.annot.muts.avinput.exonic_variant_function
HCC1143_vs_HCC1143_BL.annot.muts.avinput.log
HCC1143_vs_HCC1143_BL.annot.muts.avinput.variant_function
```

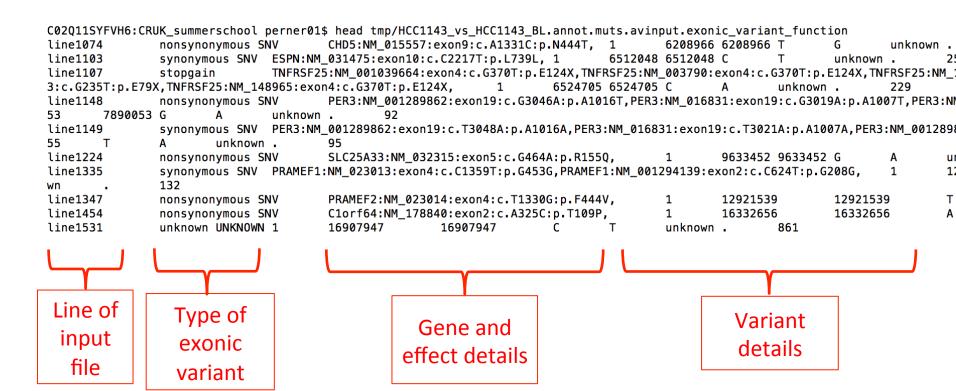
Gene-based annotation

- Output:
 - variant function file:

```
C02Q11SYFVH6:CRUK_summerschool perner01$ head tmp/HCC1143_vs_HCC1143_BL.annot.muts.avinput.variant_function
intergenic
                NONE(dist=NONE), DDX11L1(dist=1724)
                                                                                                  unknown .
                                                                  10150
                                                                          10150
                                                                                  C
                                                                                                                   449
intergenic
                                                         1
                                                                                          C
                NONE(dist=NONE), DDX11L1(dist=1694)
                                                                 10180
                                                                          10180
                                                                                                   unknown .
                                                                                                                   270
intergenic
                NONE(dist=NONE), DDX11L1(dist=1633)
                                                         1
                                                                 10241
                                                                          10241
                                                                                                  unknown .
                                                                                                                   172
intergenic
                                                         1
                NONE(dist=NONE), DDX11L1(dist=1583)
                                                                 10291
                                                                          10291
                                                                                                                   191
                                                                                                   unknown .
intergenic
                NONE(dist=NONE), DDX11L1(dist=1559)
                                                         1
                                                                 10315
                                                                          10315
                                                                                                                   251
                                                                                                  unknown .
intergenic
                NONE(dist=NONE), DDX11L1(dist=1526)
                                                         1
                                                                 10348
                                                                          10348
                                                                                                                   395
                                                                                                  unknown .
intergenic
                NONE(dist=NONE), DDX11L1(dist=1520)
                                                         1
                                                                          10354
                                                                 10354
                                                                                                  unknown .
                                                                                                                   574
intergenic
                NONE(dist=NONE), DDX11L1(dist=1517)
                                                         1
                                                                 10357
                                                                          10357
                                                                                                  unknown .
                                                                                                                   560
                                                         1
intergenic
                NONE(dist=NONE), DDX11L1(dist=1480)
                                                                 10394
                                                                          10394
                                                                                                  unknown .
                                                                                                                   624
                NONE(dist=NONE), DDX11L1(dist=1434)
intergenic
                                                         1
                                                                 10440
                                                                          10440
                                                                                                                   600
                                                                                                   unknown .
                      Distant to closest
Variant
                                                                                   Variant
                            gene or
regions
                                                                                    details
                      overlapping gene
```

Gene-based annotation

- Output:
 - exonic_variant_function file:



Exercises

- Check how many variants/what percentage of variants fall in intergenic or exonic regions?
- What is the most common exonic variant type?
- Which variants affect your favourite gene (e.g. TP53)?

Region-based annotation

 Uses same script but we need to set two more parameters:

```
$ software/annovar/annotate_variation.pl \
    -regionanno \
    -build hg19 \
    -dbtype region_dbname \
    path_to_your_annovar_file.avinput \
    path_to_your_db_folder;
```

- Options for region databases, are for example:
 - cytoband, wgRna, phastConsElements46way, tfbsConsSites, gwasCatalog, genomicSuperDups
 - See also: http://annovar.openbioinformatics.org/en/latest/user-guide/region/

Exercises

- Is there a transcription factor whose binding sites are often hit by mutations?
- Has any of the variants been found as being associated with the cancer?
- How many variant should we treat we caution because they fall into segmental duplications?

Filter-based annotation

 Uses same script but we need to change two parameters:

```
$ software/annovar/annotate_variation.pl \
    -filter \
    -build hg19 \
    -dbtype filter_dbname \
     path_to_your_annovar_file.avinput \
     path_to_your_db_folder;
```

Output:

```
C02Q11SYFVH6:CRUK_summerschool perner01$ ls tmp/*snp*
tmp/HCC1143_vs_HCC1143_BL.annot.muts.avinput.hg19_snp138_dropped
tmp/HCC1143_vs_HCC1143_BL.annot.muts.avinput.hg19_snp138_filtered
```

Exercises

- How many SNVs would you filter based on dbSNP?
- How many based on Cosmic?

All at once

```
$ software/annovar/table annovar.pl \
   path to your annovar file.avinput \
   -dbtype region dbname \
   -buildver hg19 \
   -out path to outfile.annovar \
    -remove \
    -protocol refGene, cytoBand, gwasCatalog,
genomicSuperDups,dgvMerged,snp129,esp6500si all,
cosmic70,nci60,ljb23 sift \
    -operation g,r,r,r,f,f,f,f,f \
    -nastring NA \
    -csvout;
```