# Instruction for cdot podot generation

work directory: /isilon/R_and_D/user_folders/sjung/project/oncomine

input file: 🗂 Oncomine_Myeloid_v2.20230321.GX5.hotspots.bed

manually create test.bed file: 🗂 test.bed

Run work.ipynb

```
1  import pandas as pd
2  df=pd.read_csv('test.bed',sep='\t')
3  df.fillna('-',inplace=True)
4  df['start']=df['start'].astype(int)+1
5  df['ID3']=df['ref'].astype(str)+"/"+df['alt'].astype(str)
6  df['ID4']=1
7  df[['chr','start','stop','ID3','ID4']].to_csv('vep_input.txt',index=False,sep='\t',header=False)
```

prepare an input file for VEP: 📄 vep_input.txt

e! Variant Effect Predictor - Homo_sapiens - GRCh37 Archive browser 112

Run VEP twice with different options

1. Ensembl/GENCODE transcripts and HGVS only
   a. output: ensembl.txt 📄 ensembl.txt

2. RefSeq transcripts and HGVS only
   a. output: refseq.txt 📄 refseq.txt

prepare neo transcript file (received from Sage): 📄 neo_transcript.txt

intersect the VEP outputs and the neo transcript file

```
1   # intersect between vep files and our transcripts
2   import pandas as pd
3   ref=pd.read_csv('neo_transcript.txt',sep='\t')
4   df=pd.read_csv('refseq.txt',sep='\t') #,usecols=['Location','SYMBOL','Feature','HGVSc','HGVSp',
    'Feature_original',])
5   df['Feature2']=df['Feature'].str.split('.').str[0]
6   df_merged=pd.merge(ref,df,how='inner',on=['SYMBOL','Feature2'])
7   df2=pd.read_csv('ensembl.txt',sep='\t')#,usecols=['Location','SYMBOL','Feature','HGVSc','HGVSp'])
8   df2['Feature2']=df2['Feature'].str.split('.').str[0]
9   df2_merged=pd.merge(ref,df2,how='inner',on=['SYMBOL','Feature2'])
10  # merge refseq and ensembl
11  df_concat=pd.concat([df_merged,df2_merged])
12  df_concat.to_csv('refseq_ensembl_combined.txt',sep='\t',index=False)
```

output: `refseq_ensembl_combined.txt`

```
1  #manually copy over 4th column from check.txt to test.bed for now
2  ref=pd.read_csv('test.bed',sep='\t')
3  x =pd.read_csv('refseq_ensembl_combined.txt',sep='\t',usecols=
   ['Location','SYMBOL','Feature','UPLOADED_ALLELE','HGVSc','HGVSp'])
4  x['chr']=x['Location'].str.split(':').str[0]
5  x['temp']=x['Location'].str.split(':').str[1]
6  x['start']=x['temp'].str.split('-').str[0].astype(int) - 1
7  x['start']=x['start'].astype(str)
8  x['stop']=x['temp'].str.split('-').str[1]
```

```
 9  x.to_csv('temp2.txt',sep='\t',index=False)
10  x =pd.read_csv('temp2.txt',sep='\t')
11  x['chr']=x['chr'].astype(str)
12
13  xx=pd.merge(ref,x,how='left',on=['chr','stop','UPLOADED_ALLELE'])
14  xx.to_csv('inspect2.txt',sep='\t')
15  xx.drop_duplicates(inplace=True)
16  xx.to_csv('final.txt',sep='\t')
```