

# Prioritizing Non coding Variants from Whole genome sequencing project

Sailendra Pradhananga  
DNA Club  
22/02/2017

# Drug induced myelosuppressive toxicity in Lung cancer

Lung cancer → Common cancer with high mortality rate



- Serious side effect of chemotherapy
- Multitude of genes involved in the traits

- Myelosuppression toxicity
- decrease in blood cells

**AIM : Explain variability in myelosuppression response in Lung cancer patients**

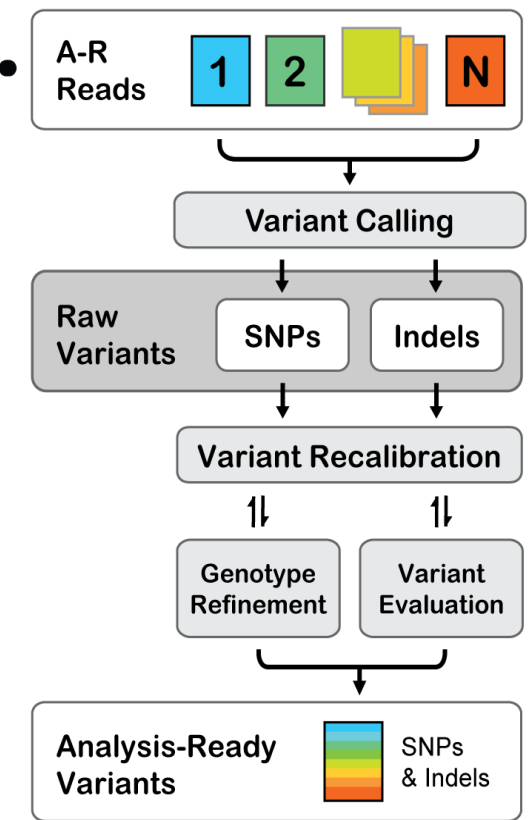
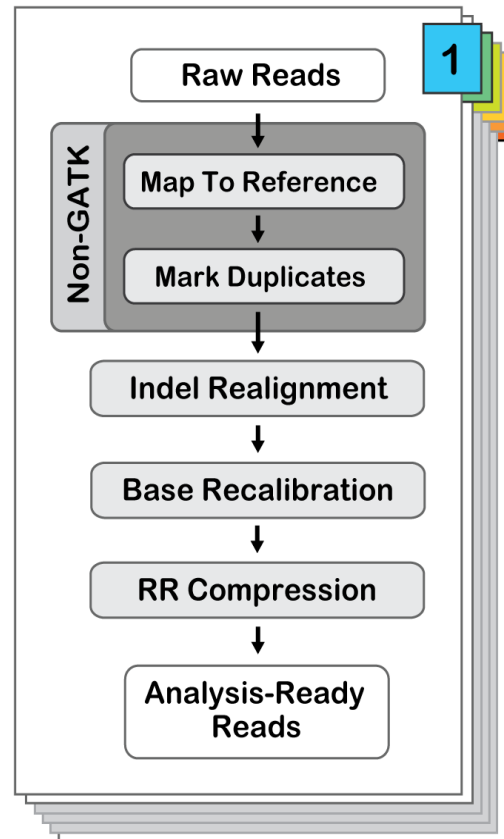
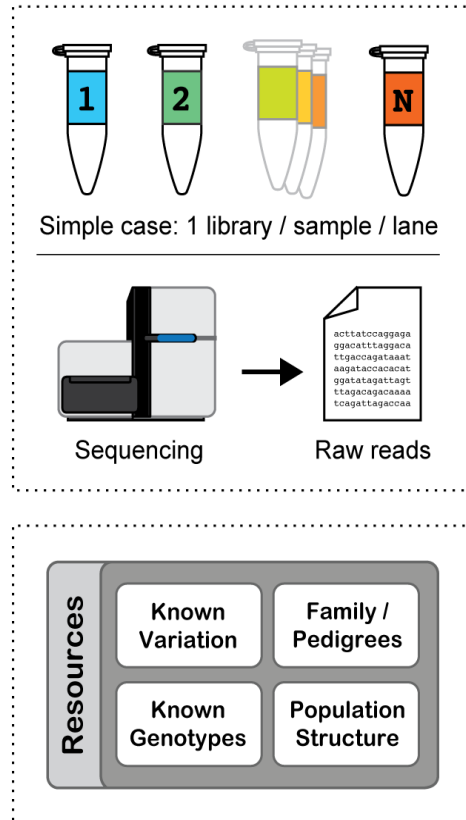
- Whole exome sequencing was performed in 215 lung cancer patients
- Few gene centric SNV variants reported and validated
- Subsequently, scaled up to whole genome sequencing of 96 extreme phenotypes of high and low toxicity

- Whole genome sequencing genotypes large number of variants in either population or individual level
- Conventional association suffers from multiple testing correction for complex traits
- Thus prioritizing of variants are done based on specific regions of genome
- In our studies, we are aiming for prioritizing low effect non coding variants for Hi-CAP probe sets

- Part1 – Comparing genetic variants from WGS 96 Lung cancer patients sample to SweGen 1000 population sequencing
- Part 2 – Annotating lung cancer non coding variants from ENCODE, FANTOM (Ref) database
- Part 3 – Classify high and low toxicity in lung cancer patients using clustering methods
- Part 4 – Enrichment and prioritizing of non coding variants in high and low toxicity phenotype groups

- Part1 – Comparing genetic variants from WGS 96 Lung cancer patients sample to SweGen 1000 population sequencing

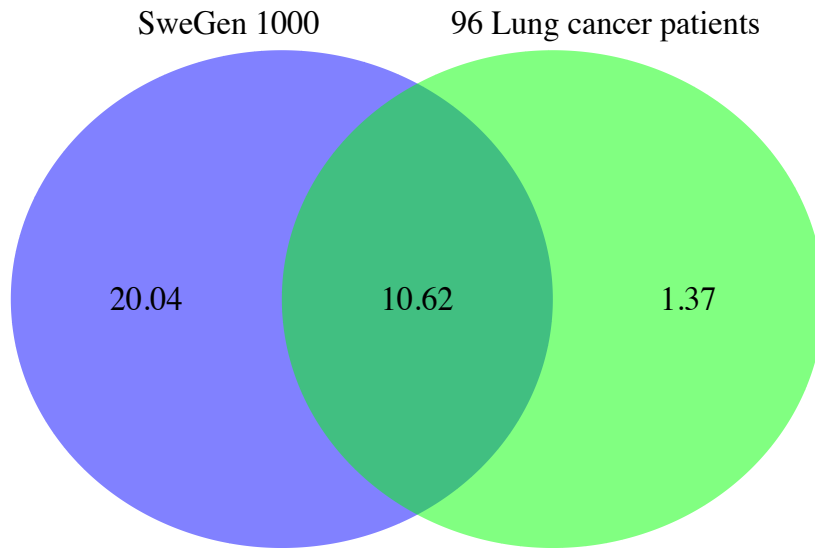
# Whole genome sequencing pipeline



Raw final Variant call in all samples

# Comparison of SweGen 1000 with 96 lung cancer sample

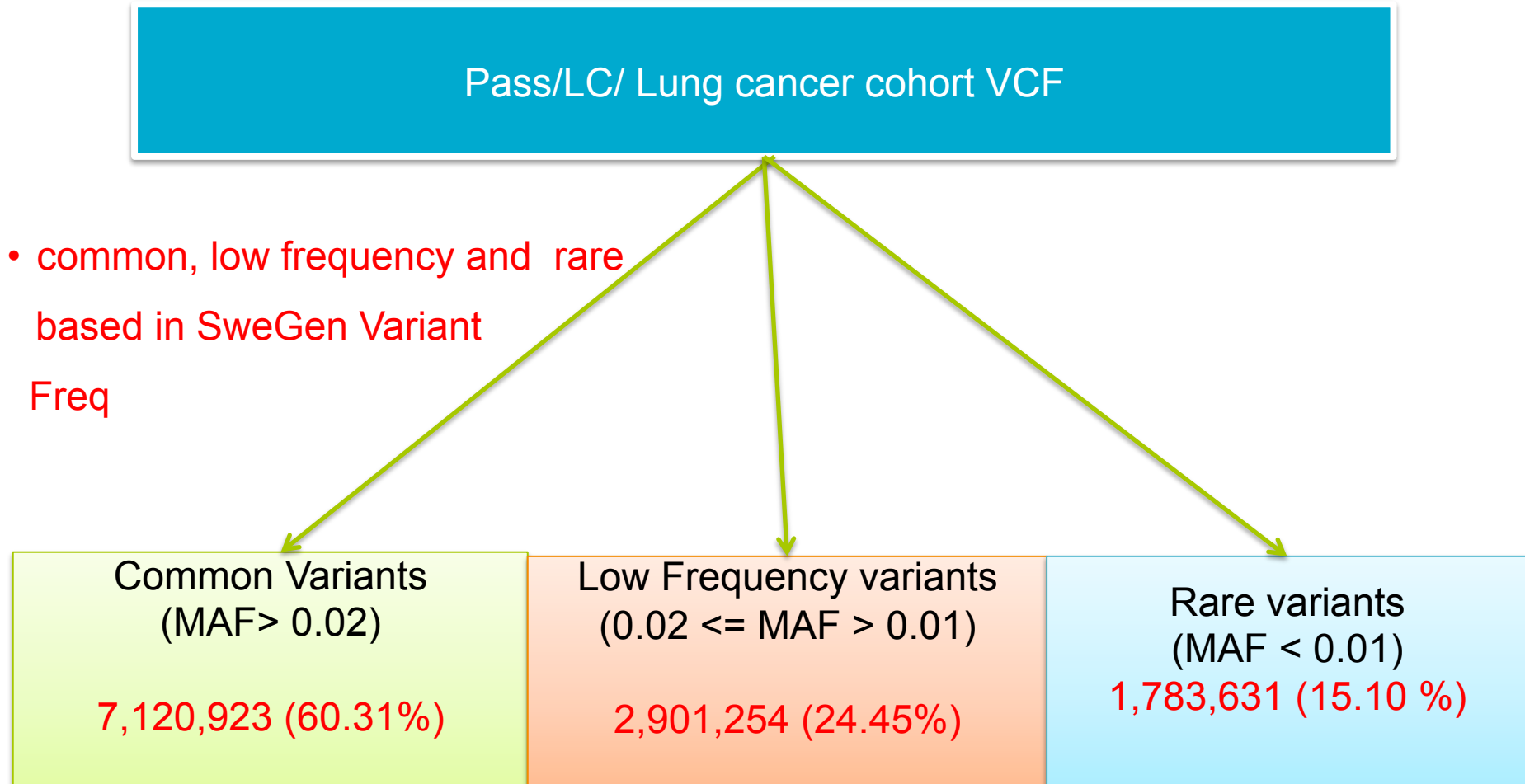
---



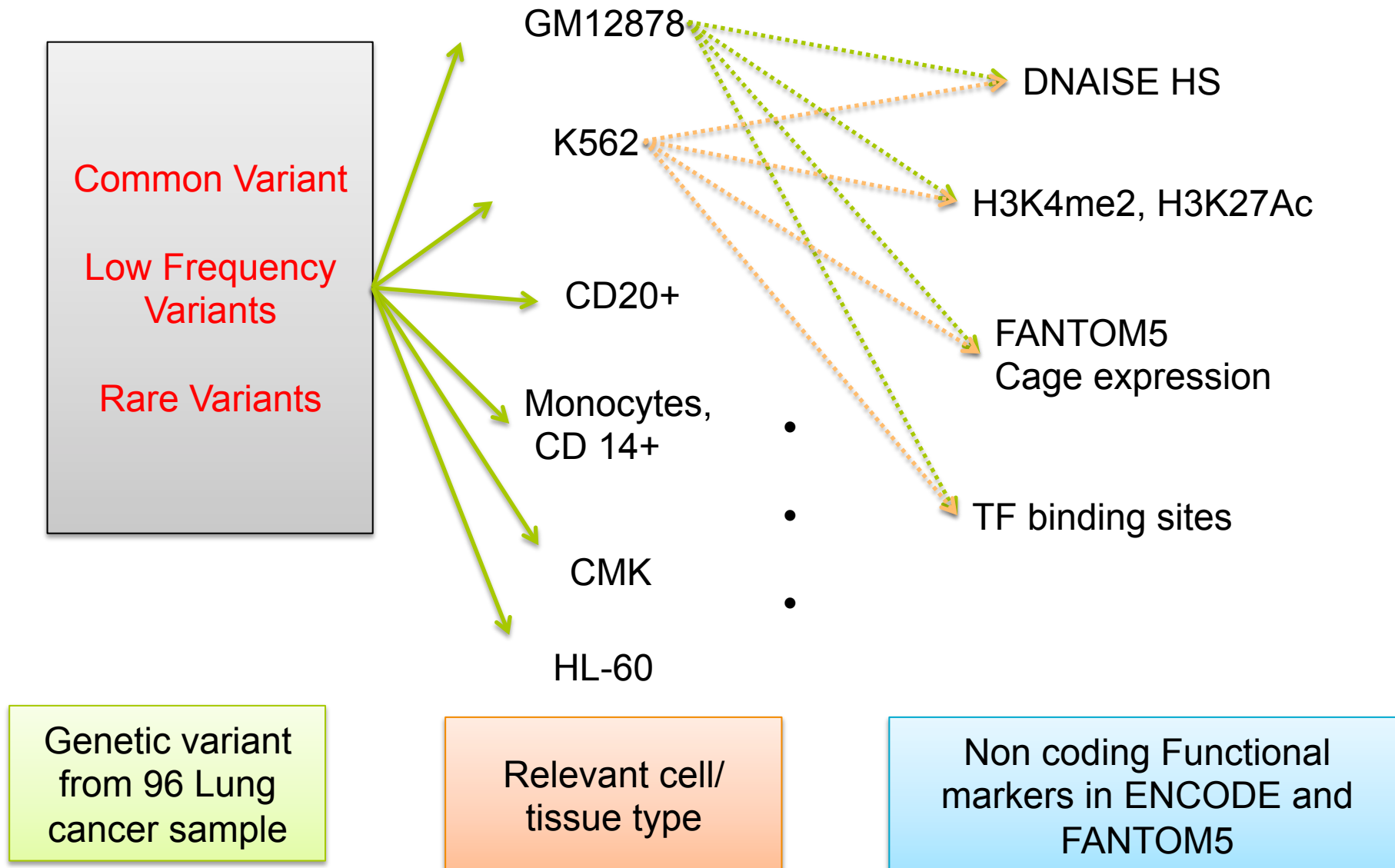
SNV counts bet two datasets in Million

- **70,353 (7.35 % ) of novel variants shared by at least 2 individuals**
- **Imputed Allele frequency from SweGen to Lung cancer cohort**

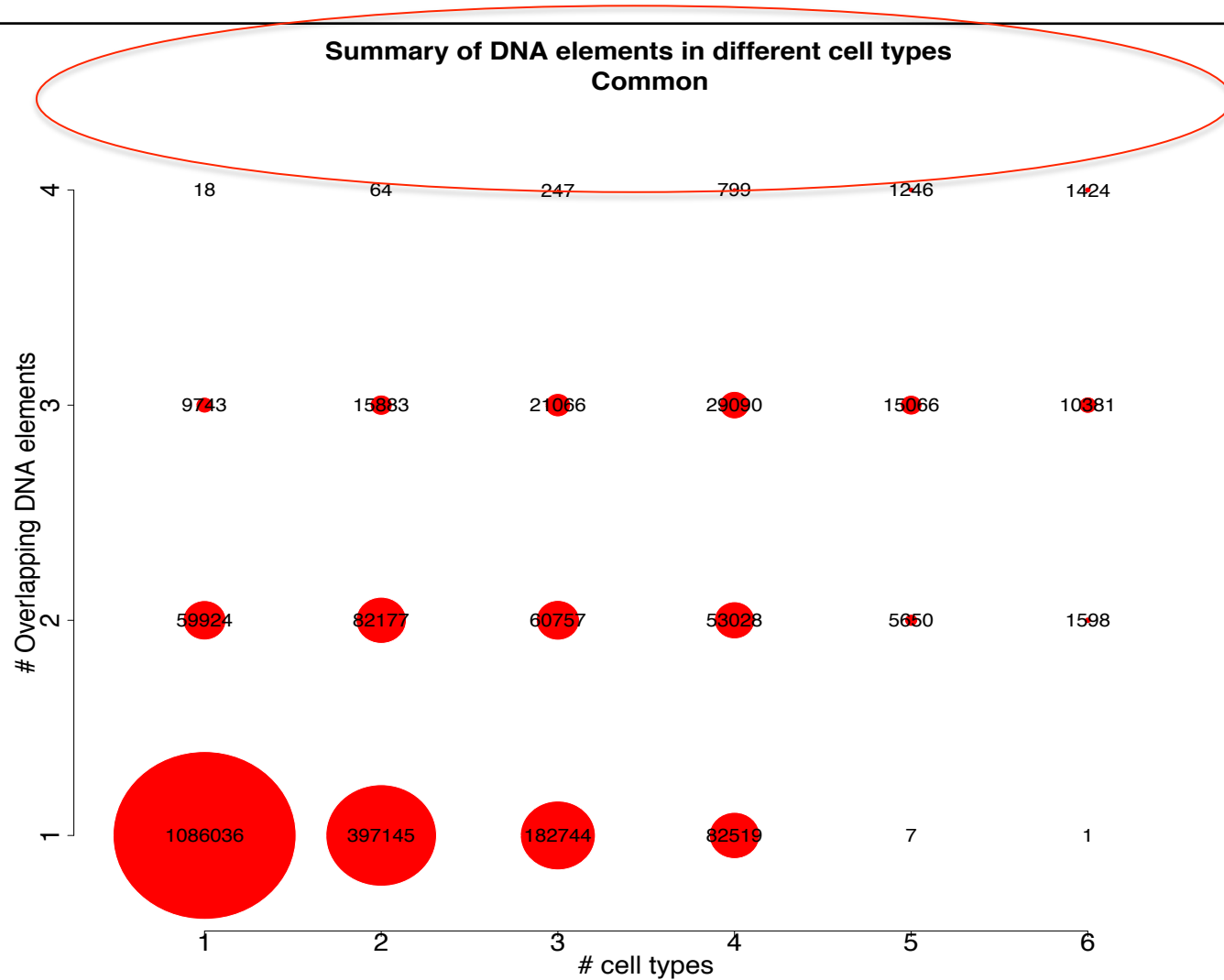




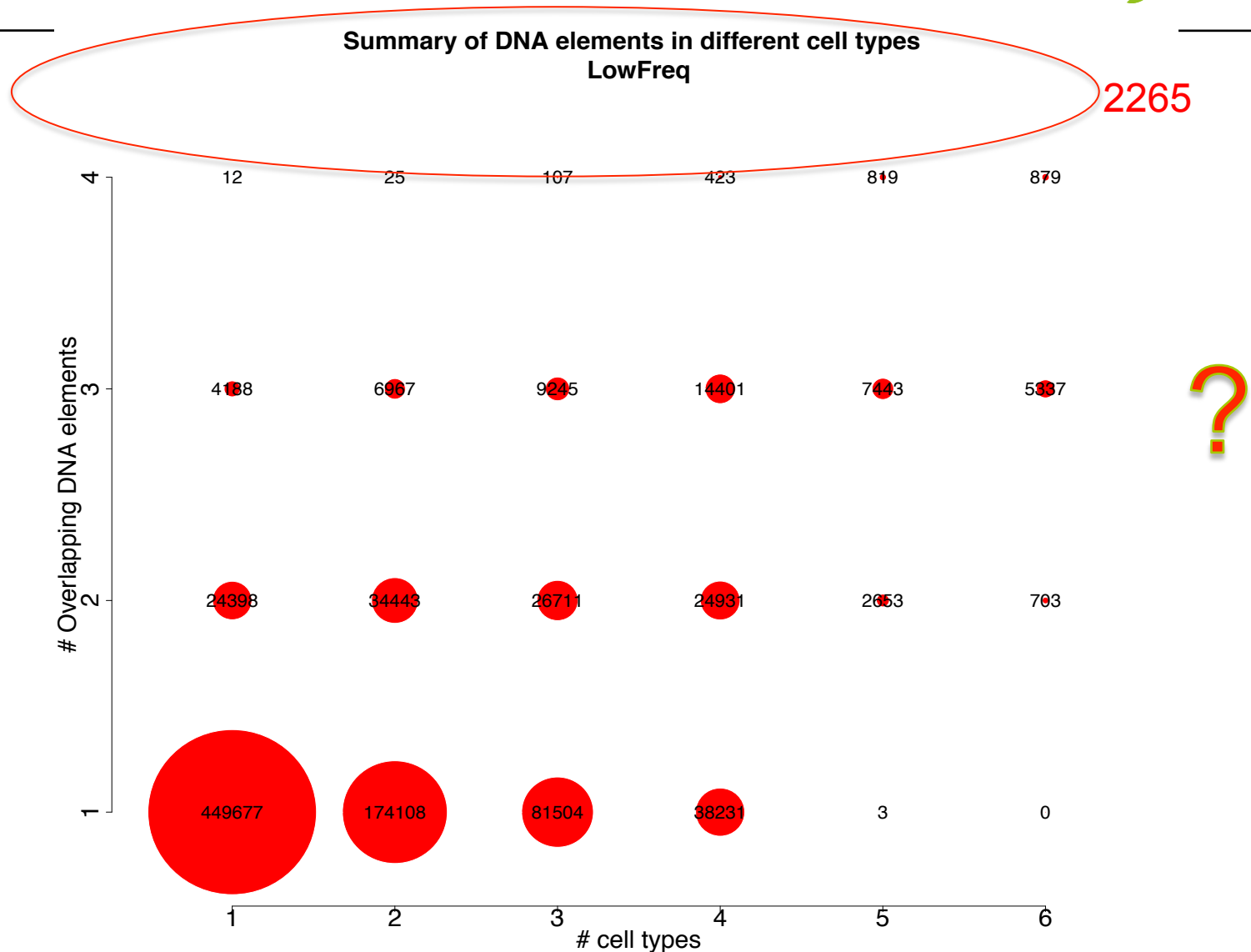
- Part 2 – Annotating lung cancer non coding variants from ENCODE, FANTOM (Ref) database

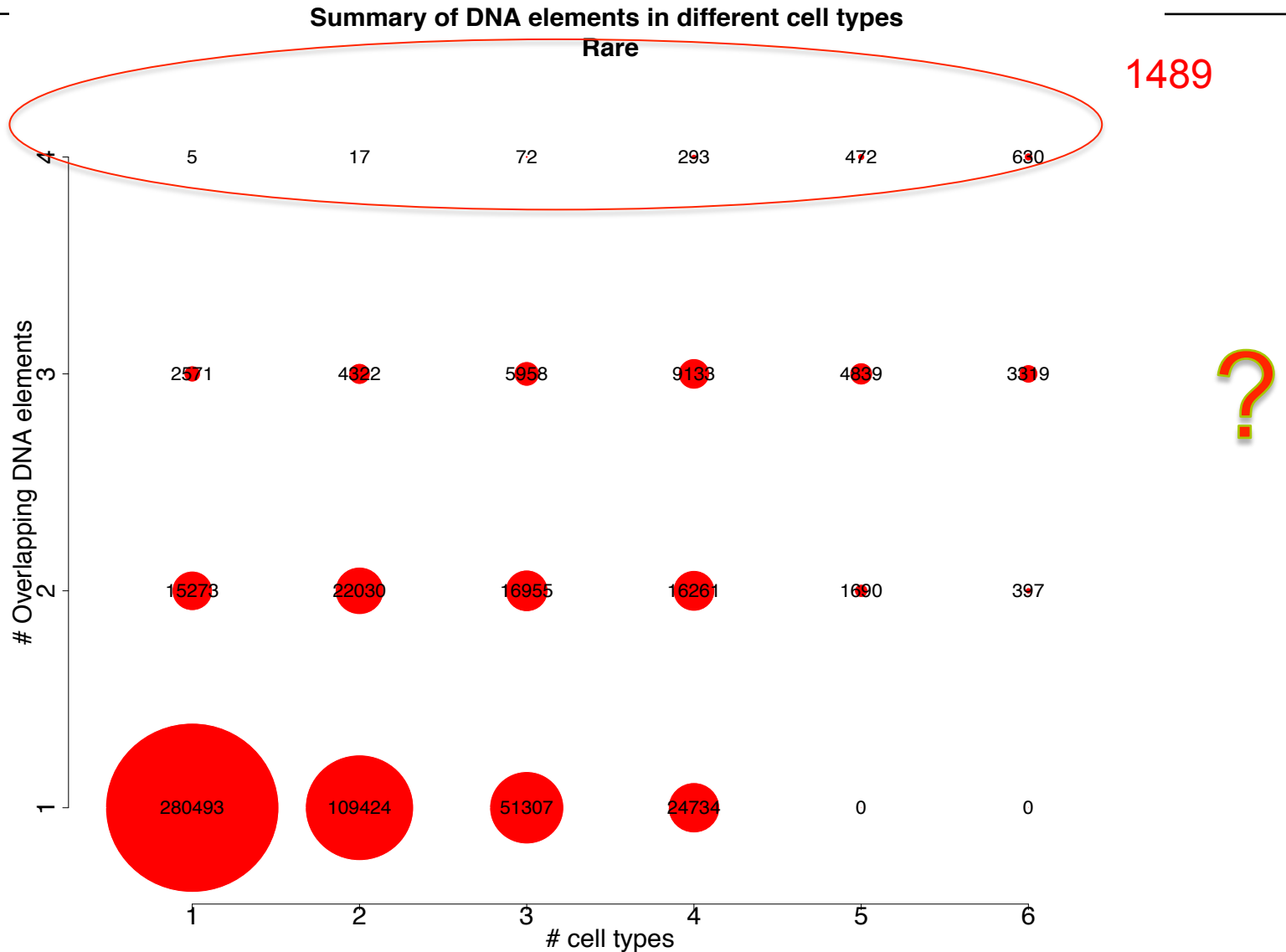


- For each variants we annotated with different functional noncoding markers of relevant cell types.
- In the initial study we annotated variant with individual cell type.
- Subsequently, aggregated approach to find variants active in different cell-types with different functional elements.



2,116,613 Common variants annotated with non coding functional marks





570,195 rare variants annotated with histone marks

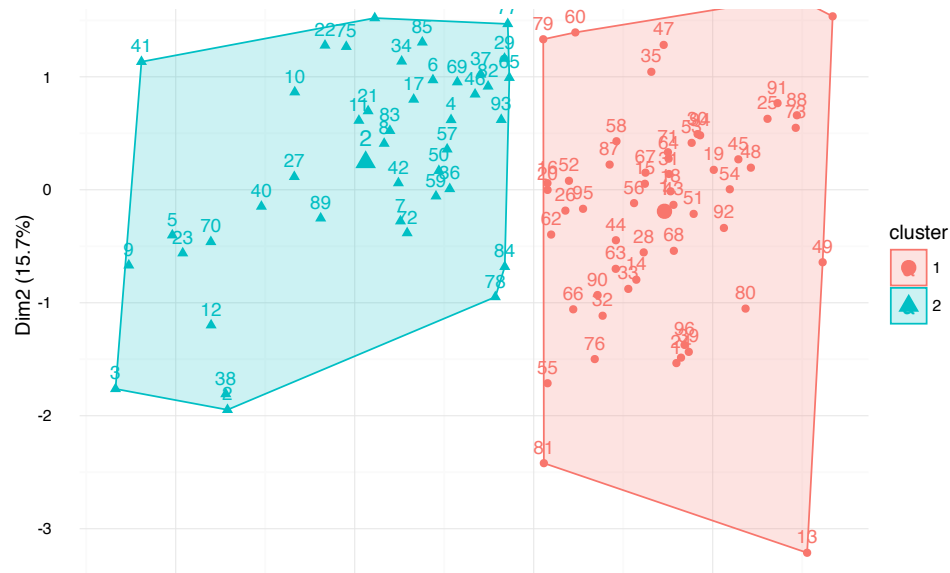
- To devise and prioritize variants based on low and high toxicity of lung cancer patients
- We used the blood cell count values **platelets(TPK)**, **leucocytes(LPK)**, **neutrophils(NPK)** after drug administration
- Cluster each variants in high and low toxicity for each **TPK**, **LPK** and **NPK** phenotype and **combined** phenotype



- Part 3 – Classify high and low toxicity in lung cancer patients using clustering methods

# Clustering results

- Using unsupervised Kmeans clustering at K=2
- Blood count values of all phenotype as an input



HT

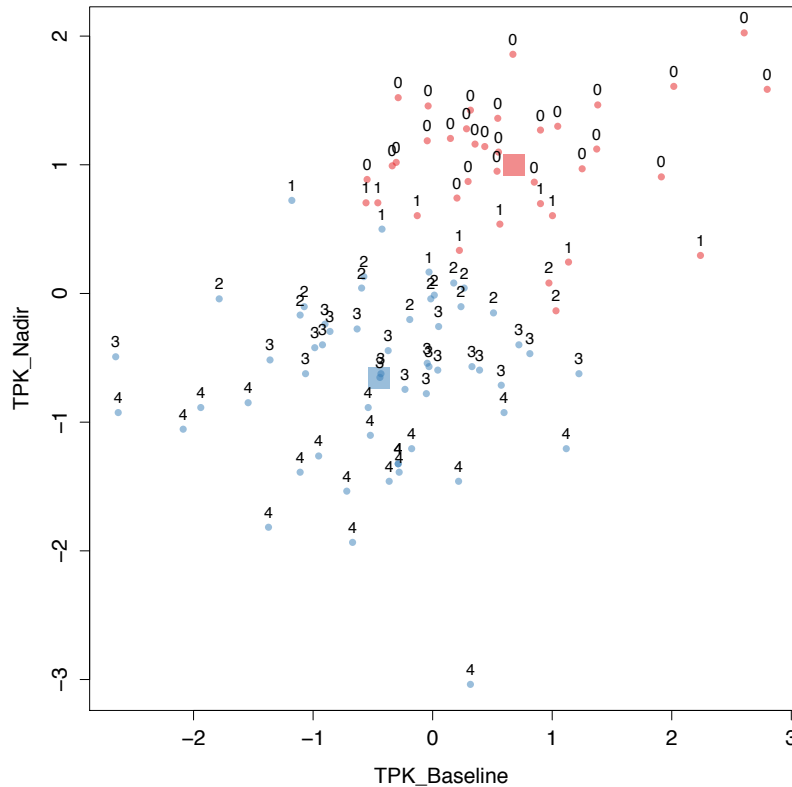
54

LT

42

# For the individual phenotypes

Clustering on 2 clustering in TPK phenotype



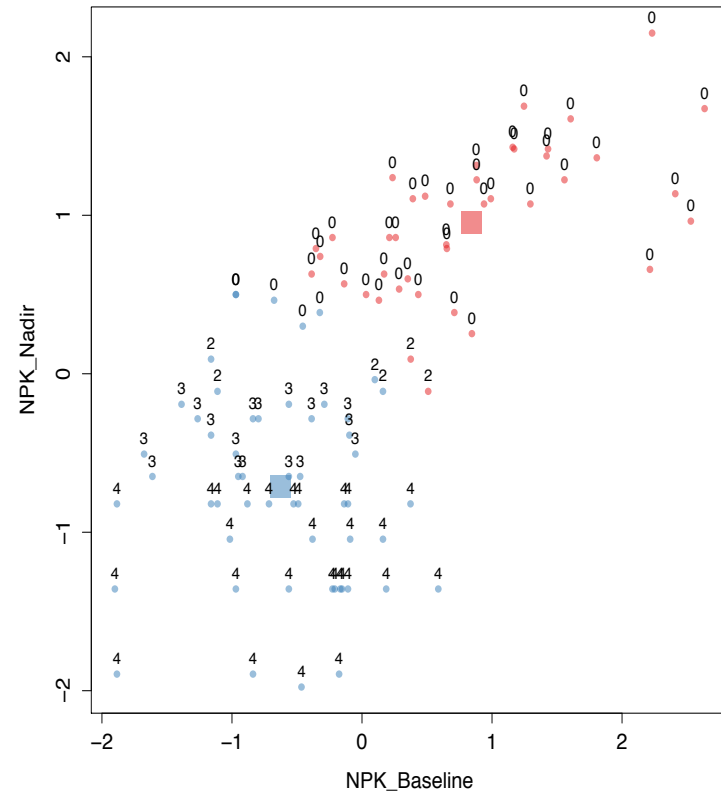
LT (red)

38

HT (Blue)

58

Clustering on 2 clustering in NPK phenotype



HT(Blue)

55

LT(Red)

41

- Part 4 – Enrichment and prioritizing of non coding variants in high and low toxicity phenotype groups

# Imputing toxicity score in each phenotypic groups

Hypothesis:

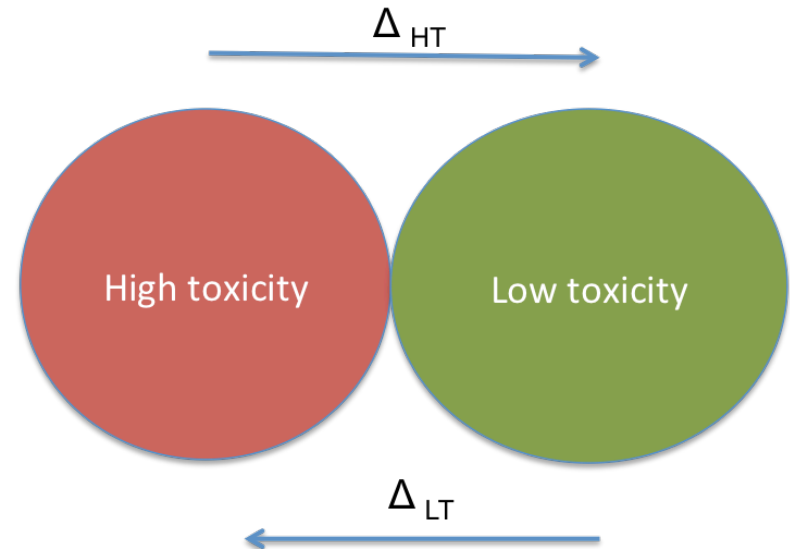
In each toxicity group enrichment of variants is attributed in delta value which is ratio of minor allele frequency in each group.

- For each variants in all group

Delta score defined as:

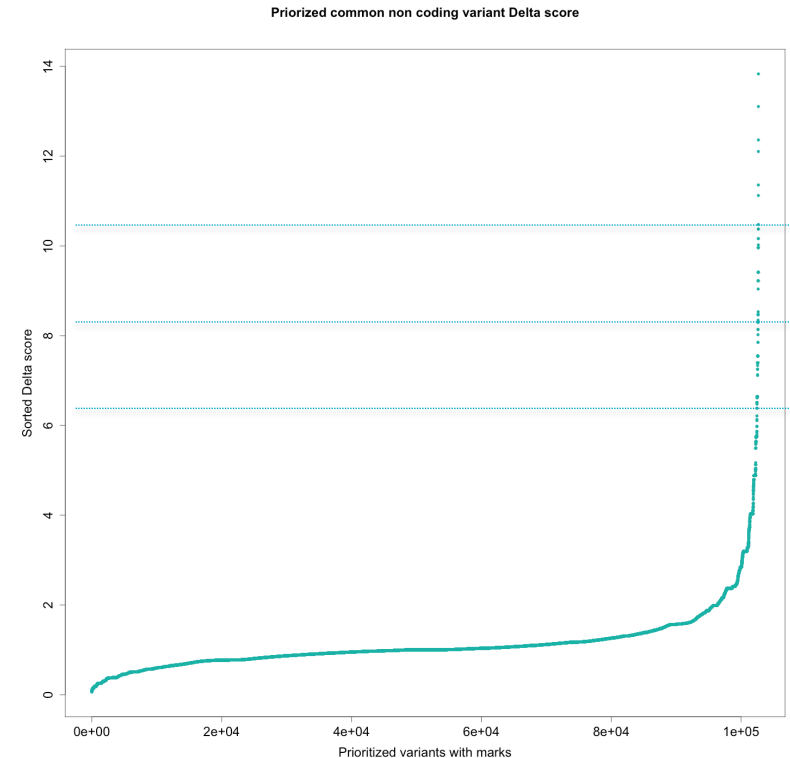
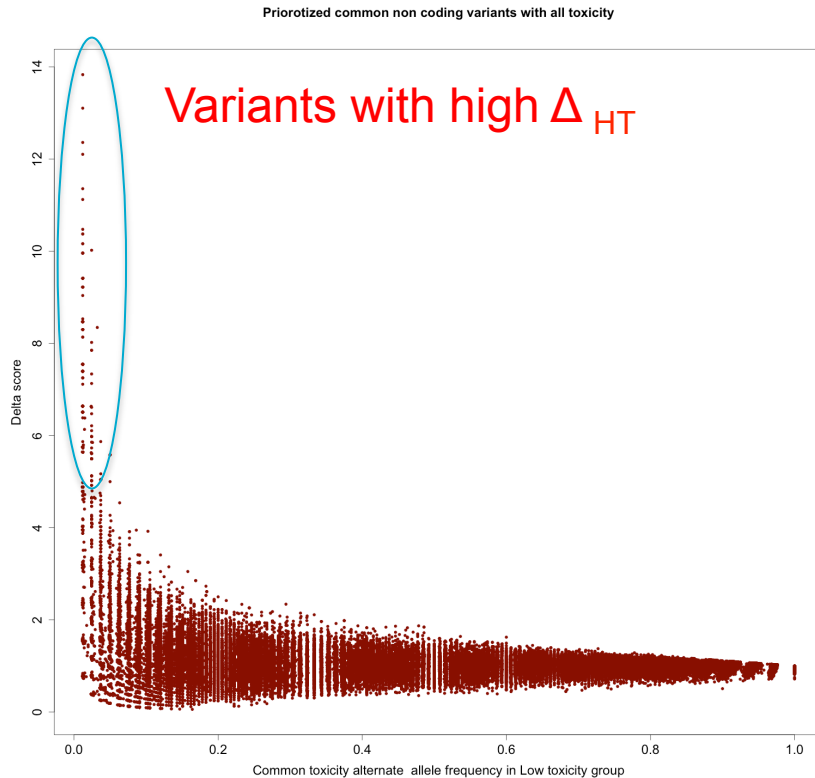
$$\Delta_{HT} = (MAF_{HT}/MAF_{LT})$$

$$\Delta_{LT} = (MAF_{LT}/MAF_{HT})$$



# Toy example

## Common ann\_level 3 / phenotype all



$\Delta_{HT}$  as high as 14 observed in Common variant levels

However still thresholding of Delta is not certain

- Similar pattern observed in all other phenotypes
- Enrichment of variants were observed in two groups

We already had variants that have been annotated with different functional markers in different cell types. So basic premise is that with increase of annotation level in variants, we are adding to functionality to variants with removal of random variants. In order to test the hypothesis we define the following terms a and b in which is defined as :

$$a = \frac{\# \text{ of non - annotated variants with } \Delta > \text{threshold}}{\text{Total number of variants}}$$

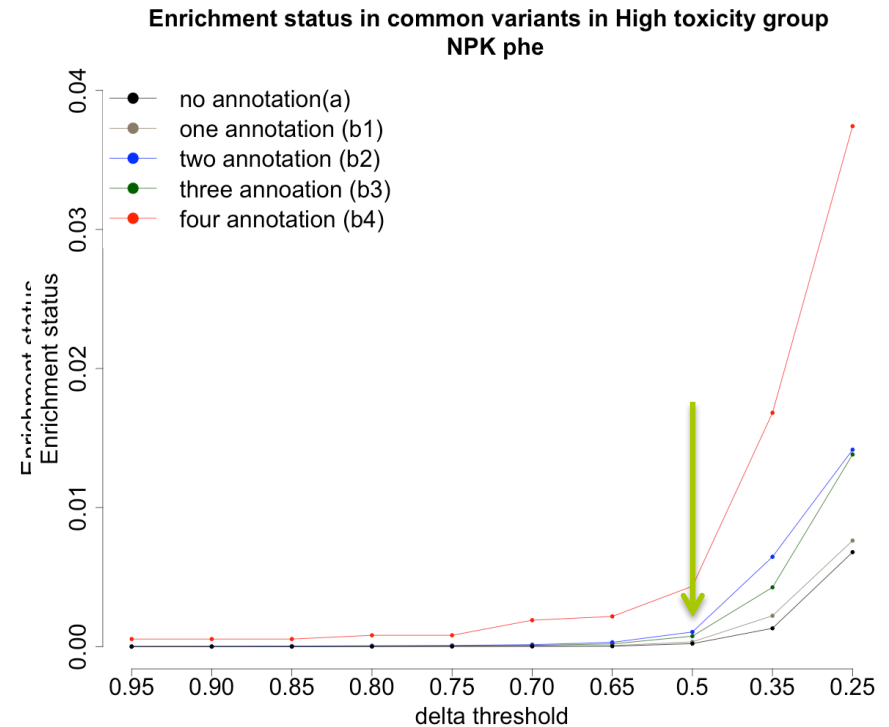
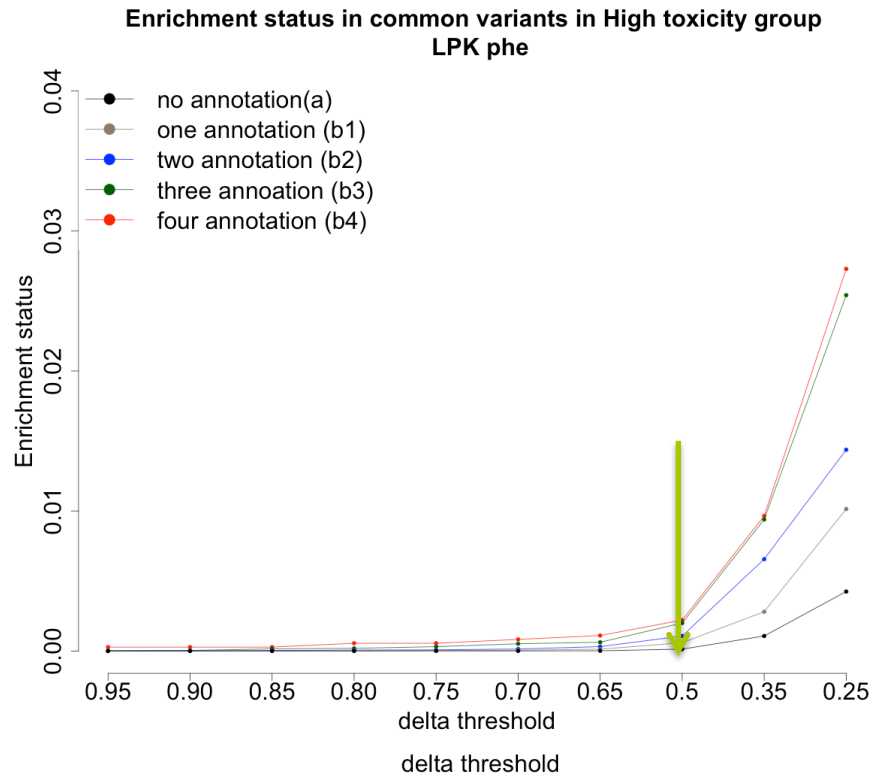
$$b = \frac{\# \text{ of annotated variants with } x \text{ levels \&\& } \Delta > \text{threshold}}{\text{Total number of variants with } x \text{ levels}}$$

where  $x = 1, 2, 3, 4$

where  $\text{threshold} = 0.95, 0.90, 0.85, 0.80, 0.75, 0.70, 0.65, 0.5, 0.35, 0.25$

Enrichment status  $b/a$  (enrichment status) in different phenotype





Trading off between number of variant and annotation level

Delta Threshold of 0.5 and Annotation level = 3

- Using threshold of 0.5 and Annotation level of 3, we prioritized **350** variants with functionality marks for high toxicity from **2million common variants**
- Similarly, using the same principal we have prioritized **293 enriched** variants with non coding function for all low toxicity groups

- Prioritized variant list of 3000 from 15 M variants
- Design probe-set for Hi-CAP studies with prioritized variants for selected cell lines MolM1, CMK
- Study interaction dynamics in variants with HiCAP analysis

- Prof. Joakim Lundeberg
- Asst Prof. Pelin Sahlen
- Assoc Prof. Henrik Green
- Niclas Bjorn
- Complex disease group

