

# Variant Annotation for TOPMed

Deepti Jain

TOPMed Analysis Workshop

August 8, 2017

# Overview of variant annotation session

- Part I – Introduction to variant annotations for TOPMed
- Part II – Defining and filtering aggregation units using variant annotations
- Part III – Hands on exercise for generating variants list used for aggregation tests
- Part IV – Hands on exercise for conducting association tests using aggregate units

# Part I

## Introduction to variant annotations for TOPMed

# What is variant annotation?

- Annotation: a note of explanation or comment
  - is the variant
    - Within a gene or intergenic
    - Within which transcript
    - Conserved
    - Deleterious
    - ...

# How are annotations used?

## 1. Rare variant association tests

- To define and filter aggregation units
- Use as weights

## 2. Fine map novel and previously known significantly associated loci to identify likely causal variants

# What is the annotation source?

- Lots of resources!

Annotation can be generated by any lab or consortium

- NCBI
- Ensemble
- UCSC
- ENCODE
- Roadmap Epigenomics Consortium
- FANTOM5
- ...

# WGSA

*J Med Genet.* 2016 February ; 53(2): 111–112. doi:10.1136/jmedgenet-2015-103423.

## **WGSA: an annotation pipeline for human genome sequencing studies**

**Xiaoming Liu<sup>1,2</sup>, Simon White<sup>3</sup>, Bo Peng<sup>4</sup>, Andrew D. Johnson<sup>5,6</sup>, Jennifer A. Brody<sup>7</sup>, Alexander H. Li<sup>1</sup>, Zhuoyi Huang<sup>3</sup>, Andrew Carroll<sup>8</sup>, Peng Wei<sup>1,9</sup>, Richard Gibbs<sup>3</sup>, Robert J. Klein<sup>10</sup>, and Eric Boerwinkle<sup>1,2,3</sup>**

<https://sites.google.com/site/jpopgen/wgsa/>

# Annotation for TOPMed variants

- Generated by Xiaoming Liu
- Variants
  - Includes 236,191,939 variants combined over TOPMed freeze 2, 3 and 4 (SNVs and Indels)
  - Includes variants flagged as failed by IRC pipeline
- Annotations
  - Used WGSA v065
  - 1,502 annotations
- Typically an annotation set is generated for every IRC freeze
- All WGSA annotation releases are made available in the exchange area



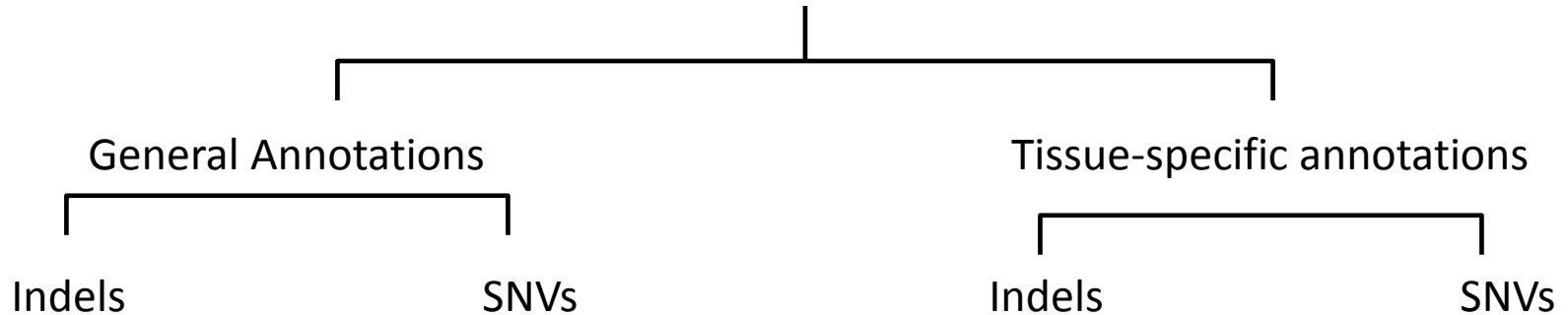
# 1,502 annotations for TOPMed

- Gene based location and consequence
  - Softwares : SnpEff, ANNOVAR, VEP
  - Gene models: Ensembl ,RefSeq ,UCSC
- Transcript-specific annotation
- Loss-of-function annotations (eg: LOFTEE)
- Deleteriousness predictions(CADD, MetaSVM, ssSNV etc)
- Allele frequencies (1000G, UK10K, EXAC etc)
- Regulatory annotations (ENCODE, Roadmap, FANTOM5)
- Conservation scores
- Mappability scores
- rsIDs
- Many more ....

**It is guaranteed that you will not use all of the annotations for an analysis.**

**We recommend starting with a subset of frequently used annotations**

# TOPMed WGS annotation directory structure



- 23 .gz file (one per chromosome)
- 1 .txt Data dictionary
- 1 .tsv file with first 1000 lines of chr1

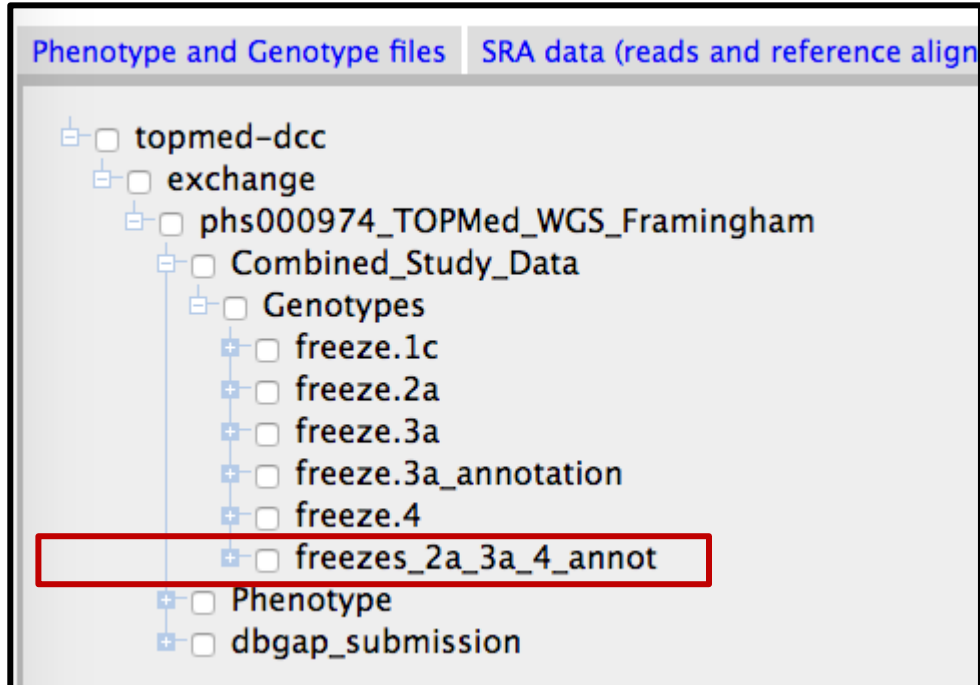
$$\times 4 = 100$$

The size of the annotation directory is **102G**  
Chr1 general annotation file for SNVs is **5GB**

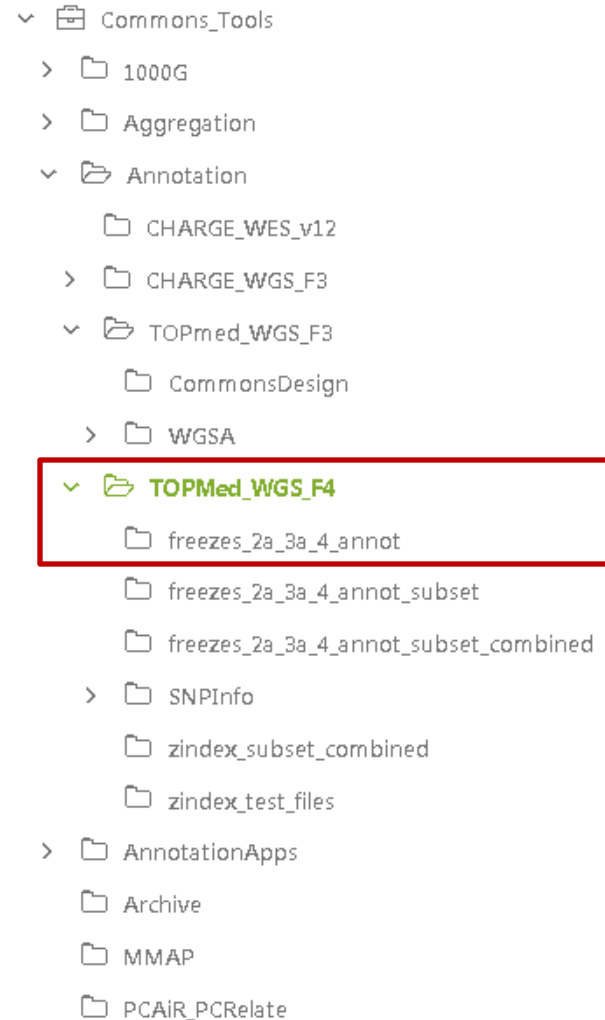
*Review data dictionary “List of resources v0.65.pdf” before using annotations*

# Where to find WGS data annotations for TOPMed?

## Exchange Area



## Analysis commons



# Gene-based annotations are at transcript level

## chr:10273 T>C

### VEP\_ensembl\_Transcript\_ID

ENST00000456328|ENST00000488147|ENST00000438504|ENST00000515242|ENST00000541675|ENST00000423562|ENST00000450305|ENST00000538476|ENST00000518655

### VEP\_ensembl\_Consequence

upstream\_gene\_variant|downstream\_gene\_variant|downstream\_gene\_variant|upstream\_gene\_variant|downstream\_gene\_variant|downstream\_gene\_variant|upstream\_gene\_variant|downstream\_gene\_variant|splice\_region\_variant

### VEP\_ensembl\_Gene\_Name

DDX11L1|WASH7P|WASH7P|DDX11L1|WASH7P|WASH7P|DDX11L1|WASH7P|DDX11L1

### VEP\_ensembl\_Gene\_ID

ENSG00000223972|ENSG00000227232|ENSG00000227232|ENSG00000223972|ENSG00000227232|ENSG00000227232|ENSG00000223972|ENSG00000227232|ENSG00000223972

### Ensembl\_Regulatory\_Build\_Overviews

ctcf

### VEP\_ensembl\_LoF

.|.|.|.|.|.|.|.|HC

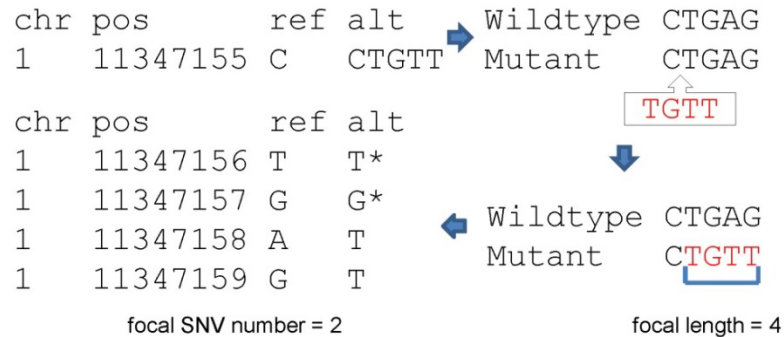
NOTE: Missingness is denoted as “.” in all annotation fields

# For some annotations indels are translated into SNVs

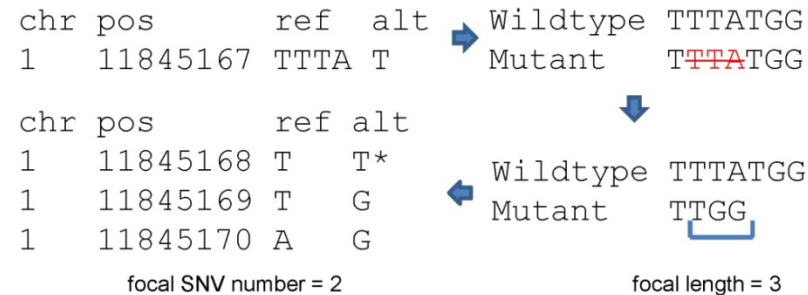
GERP++ RS score  
 0.325{1}0.097{1}0.392{1}2.010{1}

Genomic Evolutionary Rate Profiling (GERP) identifies constrained elements in multiple alignments by quantifying substitution deficits. the larger the score, the more conserved the site

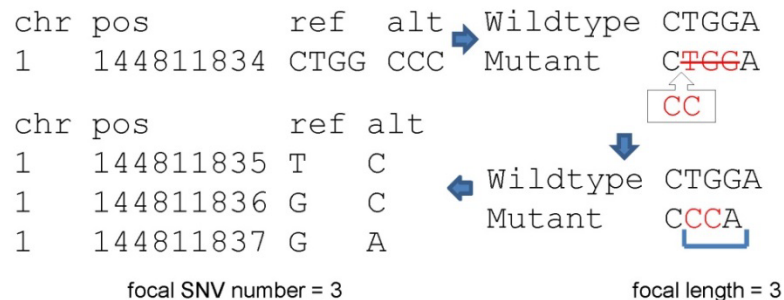
## A: insertion



## B: deletion



## C: replacement



## Part II

# Defining and filtering aggregation units using variant annotations

# How are annotations used?

## 1. Rare variant association tests

- **To define and filter aggregation units**
- Use as weights <sup>1</sup>

## 2. Fine map novel and previously known significantly associated loci to identify likely causal variants

<sup>1</sup>Morrison AC, Huang Z, Yu B, et al. Practical Approaches for Whole-Genome Sequence Analysis of Heart- and Blood-Related Traits. Am J Hum Genet. 2017;100(2):205-215.

## Steps involved in generating aggregate variant list for association testing

STEP1: Define aggregation units

- which genomic regions will be included in each unit

STEP2: Decide on filtering criteria

- which variants will be filtered within each unit

*Goal is to create list of variants in each aggregation unit which can be used in multiple variants association tests ( example Burden and SKAT tests)*

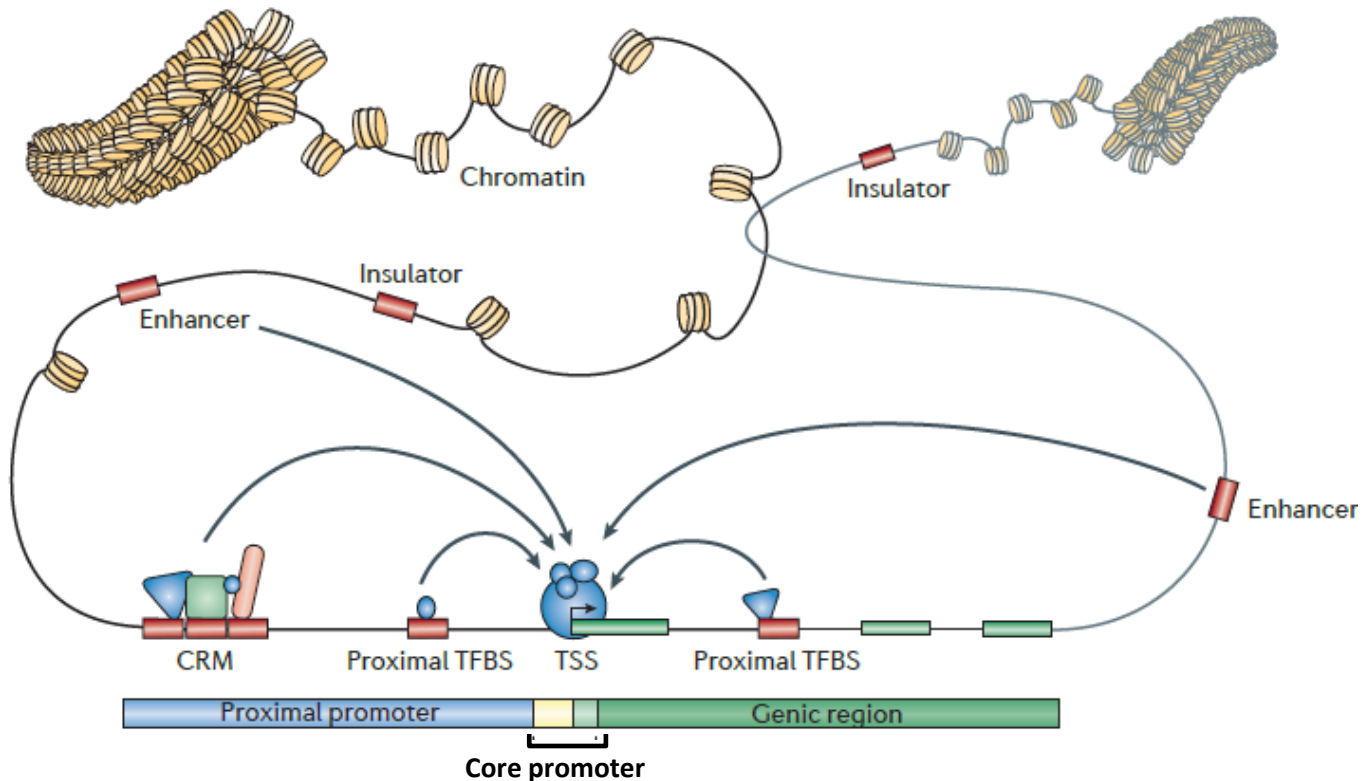


## PART II

# STEP1: Define aggregation units

Gene is one of the fundamental units of biology and gene-based aggregation units are frequently used in rare variant association testing

# What is a gene?



TFBS: transcription factor binding site,

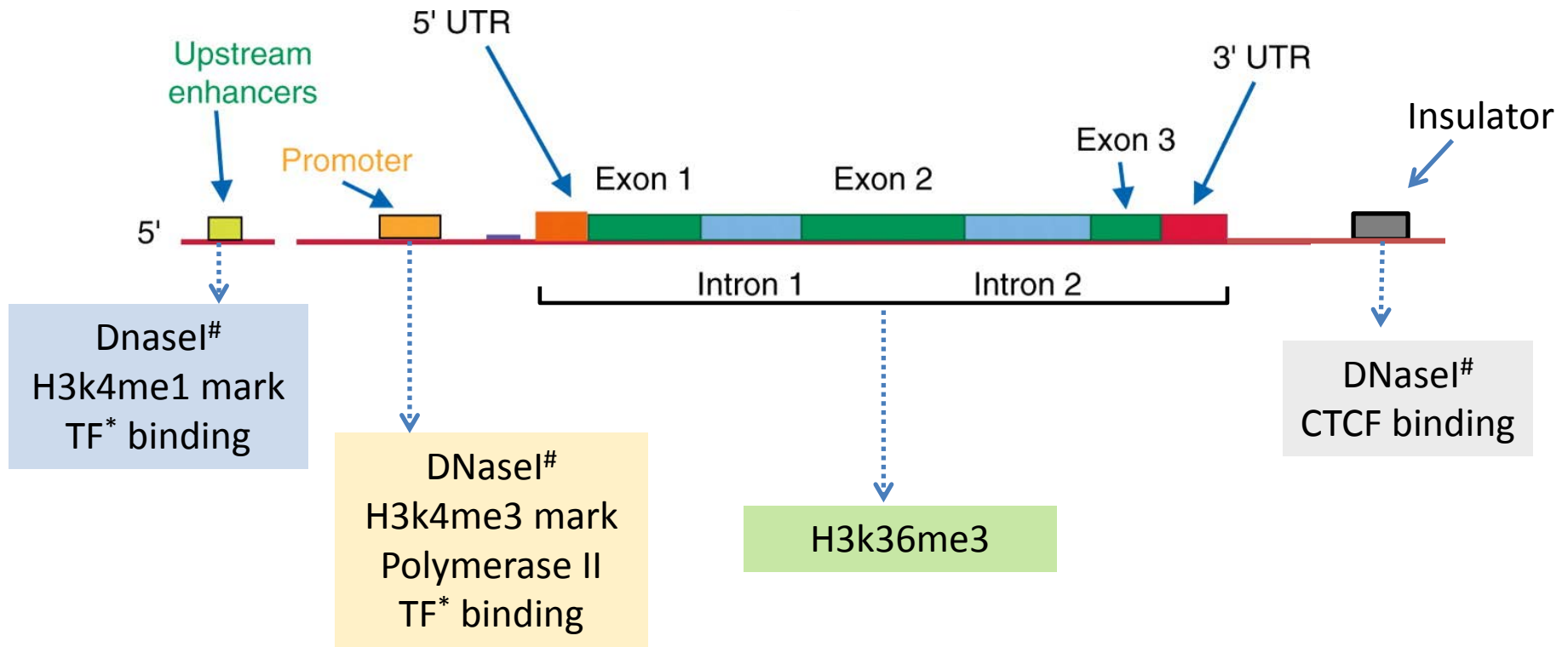
CRM: cis-regulatory module

UTR: untranslated region

Functional gene unit = transcript + promoter + enhancers

Transcript = Exon + Intron + 3'UTR + 5'UTR

# Biochemical signatures typically associated with non-coding functional elements



**Enhancer** : Interacts with promoter can be involved in repression or induction of a gene

**Promoter** : Genomic element where the transcription machinery assembles

**UTR** : Untranslated region

**EXON** : Coding part of a transcript (mRNA)

**INTRON** : Non-Coding part of a transcript (mRNA)

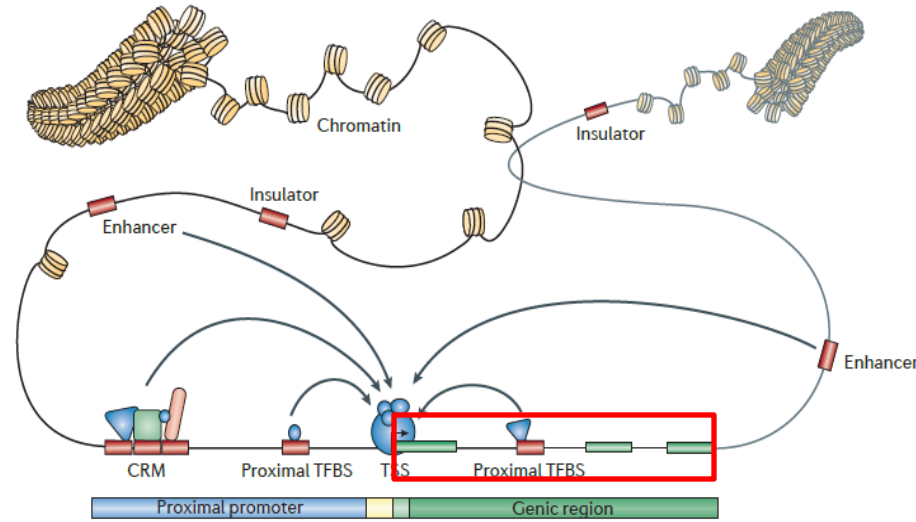
**Insulator** : Barriers that protect genes from influence of outside enhancers or inactivating chromatin structures

**NOTE: These biochemical marks are tissue-specific . Additionally, these may also show temporal and treatment specific variations within a cell type**

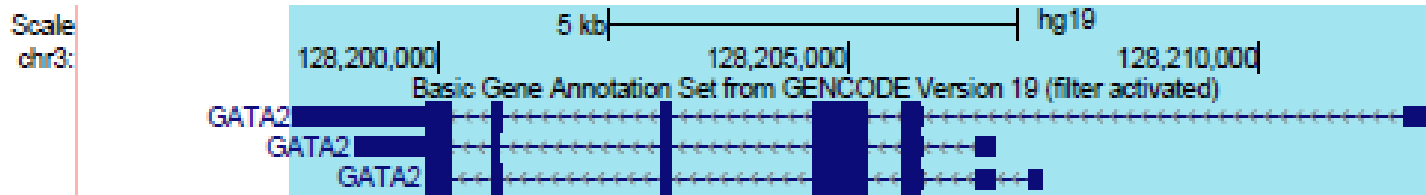
\* TF : transcription factor, # DNaseI Hypersensitivity, which is an indicator of chromatin accessibility

Functional gene unit = transcript + promoter + enhancers

# Transcript and gene models

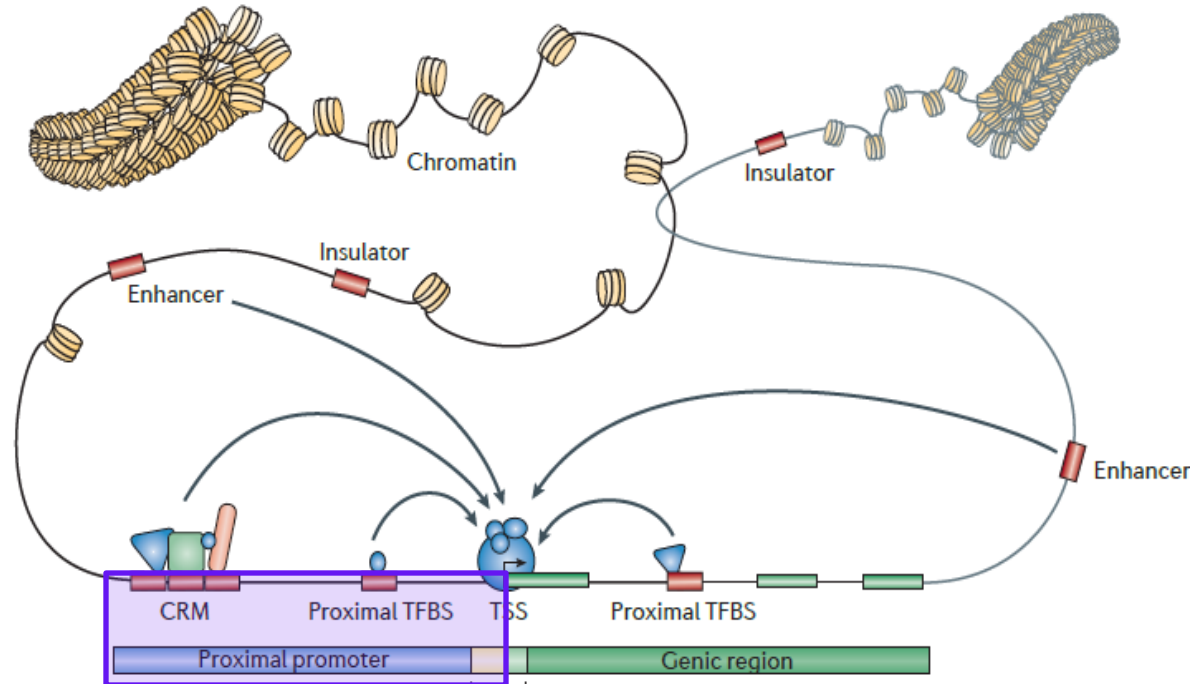


- Gene models: GENCODE, RefSeq, UCSC



- Gene
  - Genomic region spanning all the transcripts of a gene
  - A single transcript (example expressed transcripts in RNA-Seq data of relevant tissue)

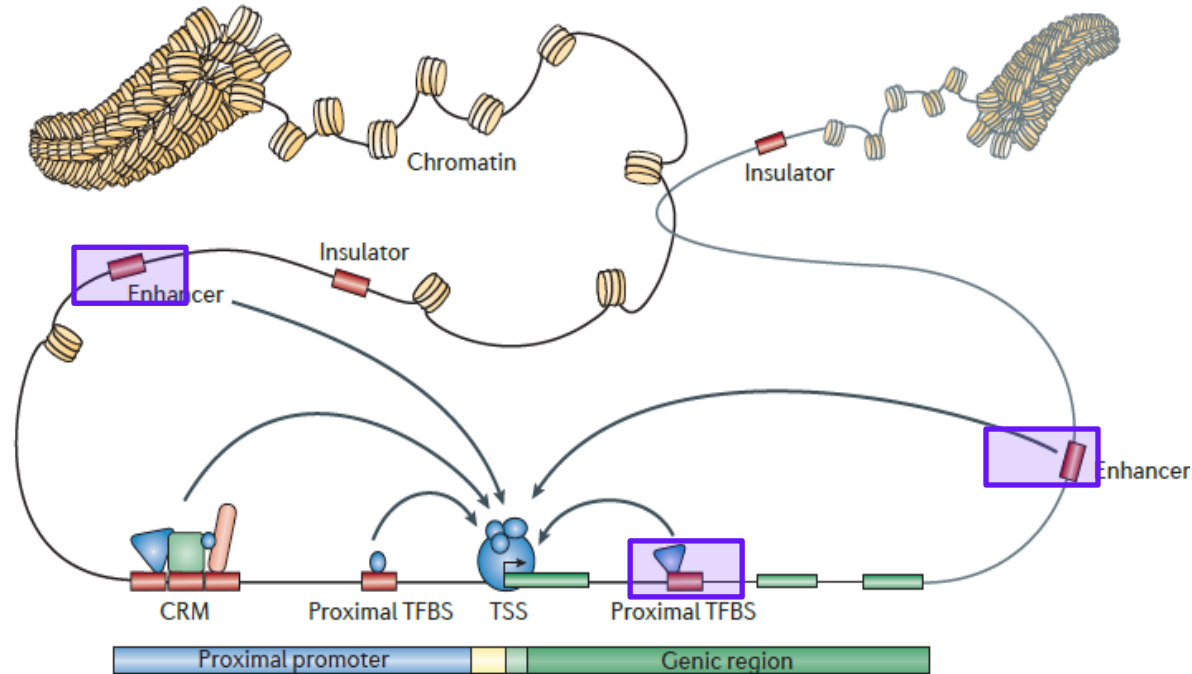
# Promoters



- Some distance upstream from TSS (typically 5Kb)
- 5Kb upstream overlapping with H3K4me3 and or H3K27ac mark
- 5Kb upstream overlaps with DNaseI hypersensitive regions
- 5Kb upstream that overlaps with CAGE peaks<sup>1</sup>

<sup>1</sup>Morrison AC, Huang Z, Yu B, et al. Practical Approaches for Whole-Genome Sequence Analysis of Heart- and Blood-Related Traits. *Am J Hum Genet.* 2017;100(2):205-215.

# Enhancers



- Flanking regions overlapping with H4K4me1 and or H3K27ac
- Flanking regions overlapping with DNaseI hypersensitive regions
- Enhancer-gene link predictions<sup>1,2</sup>
- Chromosome conformation capture (3C,4C,Hi-C etc.)

<sup>1</sup>Thurman RE et.al *Nature*. 2012 Sep 6; 489(7414):75-82.

<sup>2</sup>Forrest AR, Kawaji H, Rehli M, et al. A promoter-level mammalian expression atlas. *Nature*. 2014;507(7493):462-70.

# Example gene-based aggregation units

- Gene + flanking regions
- Gene + enhancer + promoter
- UTR's+ enhancer + promoter
- Promoter of a gene
- First intron of a transcript

Aggregation units that are continuous genomic ranges

chromosome	start	end	gene_id
chrX	99883667	99894988	ENSG00000000003.10
chrX	99839799	99854882	ENSG00000000005.5
chr20	49551404	49575092	ENSG000000000419.8
chr1	169818772	169863408	ENSG000000000457.9
chr1	169631245	169823221	ENSG000000000460.12
chr1	27938575	27961788	ENSG000000000938.8

Aggregation units that are non-continuous genomic ranges

chromosome	start	end	gene_id	type
chr1	27938575	27961788	ENSG000000000938.8	gene
chr1	27945001	27945200	ENSG000000000938.8	7Enh
chr1	27945201	27945400	ENSG000000000938.8	7Enh
chr1	27945401	27945600	ENSG000000000938.8	7Enh
chr1	27945601	27945800	ENSG000000000938.8	7Enh



# Other units of aggregation

(not covered in this workshop)

- Genes in a pathway<sup>1</sup>
- Protein-protein interaction domains<sup>2</sup>
- Topological associated domains<sup>3</sup>
- Specific histone modification marks<sup>3</sup>
- DNaseI Hypersensitive sites/nucleosome depleted regions
- ...

<sup>1</sup>Richardson TG, Timpson NJ, Campbell C, Gaunt TR. A pathway-centric approach to rare variant association analysis. *Eur J Hum Genet.* 2016;25(1):123-129.

<sup>2</sup>Richardson TG, Shihab HA, Rivas MA, et al. A Protein Domain and Family Based Approach to Rare Variant Association Analysis. *PLoS ONE.* 2016;11(4):e0153803.

<sup>3</sup>Lumley, T., Brody, J., Peloso, G., & Rice, K. (2016). Sequence kernel association tests for large sets of markers: tail probabilities for large quadratic forms. doi:10.1101/085639

## PART II

### STEP2: Decide on filtering criteria

# Scenario 1: simple filtering

- Genic unit

Transcript range + 20 kb flanking region upstream and downstream

- Filters:

CADD phred score  $\geq 10$  and MAF  $\leq 1\%$

**Combined Annotation Dependent Depletion (CADD)** is a tool for scoring the deleteriousness of variants. A scaled C-score of greater or equal 10 indicates that these are predicted to be the 10% most deleterious substitutions that you can do to the human genome

# Scenario 2: Using multi-tissue regulatory regions

- Genic unit
  - Gene + 20 kb flanking region upstream and downstream
- Filters:
  - A. Flanking region
    - Overlaps with “Ensembl\_Regulatory\_Build\_Overviews”
  - A. Gene region
    - Overlaps with “Ensembl\_Regulatory\_Build\_Overviews” OR
    - Overlaps with LOF variants
- Ensembl\_Regulatory\_Build\_Overviews
  - genome segment prediction based on 17 cell types from ENCODE and Roadmap.
  - ctcf - CTCF binding sites,
  - distal - Predicted enhancers
  - open - Unannotated open chromatin regions
  - proximal - Predicted promoter flanking regions
  - tfbs - Unannotated transcription factor binding sites
  - tss - Predicted promoters
- ENCODE\_Dnase\_cells: number of cell lines supporting a DNase I hypersensitive site

## Scenario 3: Using tissue specific regulatory regions -I

- All epigenomics marks of erythroleukemia cell line K562 (EID: E123) in WGSA
  - E123-DNase.macs2.narrowPeak
  - E123-H2A.Z.narrowPeak
  - E123-H3K27ac.narrowPeak
  - E123-H3K27me3.narrowPeak
  - E123-H3K36me3.narrowPeak
  - E123-H3K4me1.narrowPeak
  - E123-H3K4me2.narrowPeak
  - E123-H3K4me3.narrowPeak
  - E123-H3K79me2.narrowPeak
  - E123-H3K9ac.narrowPeak
  - E123-H3K9me1.narrowPeak
  - E123-H3K9me3.narrowPeak
  - E123-H4K20me1.narrowPeak

## Scenario 3: Using tissue specific regulatory regions -II

- Genic unit
  - Gene + 20 kb flanking region upstream and downstream
- Filters:
  - A. Flanking region
    - Overlaps with either H3K4me3 or H3K4me1 enriched regions ) & DNaseI hypersensitivity sites in k562 cells
  - B. Gene region
    - overlap (either H3K4me3 or H3K4me1 enriched regions ) & DNaseI hypersensitivity sites in k562 cells OR
    - Overlaps with LOF variants

## DCC will provide frequently used aggregation units on the EA

- GENCODE Gene + flanking regions
  - 0kb,5kb ,20kb,100kb,200kb
  - Filtered to keep variants overlapping “Ensembl\_Regulatory\_Build\_Overviews” and LOF variants
- Promoter
  - 5Kb upstream of GENCODE genic unit
  - Filtered to keep “tss” and “proximal” overlapping variants from “Ensembl\_Regulatory\_Build\_Overviews”
- First intron of GENCODE transcripts
- DCC will try to accommodate requests for assistance in defining units and lists of variants within them
- Requests for creating aggregation units and list of variants within them will be accepted as time permits

# Summary

- A large set of annotations are available for TOPMed through WGSA
- Annotations can be used to create and filter aggregation units
- Frequently used aggregation units will be made available by DCC on the Exchange Area



## Part III

Hands-on exercise for generating  
variants list used for aggregation  
tests

# Key libraries and functions

Parse the WGSA annotation file	
library (wgsaparsr)	Package for working with WGSA output files
get_fields	list the annotation fields available in a WGSA output file
parse_to_file	Converts list-fields into multiple rows
parse_indel_to_file	Same as “parse_to_file” but for indel annotations
Create aggregation units file using GENCODE genes	
library(genetable)	Package for working with .gtf gene model files
<b>import_gencode</b>	import the gtf file to a tidy data frame
<b>filter_gencode</b>	filter gtf file on different features and tags
<b>define_boundaries</b>	define the boundaries of the feature of interest
Aggregate variants by genic units and create input file for association testing	
<ul style="list-style-type: none"><li>- R code for the workshop</li><li>- DCC uses a MySQL server for creating and filtering variant list in aggregation units using WGSA annotations</li></ul>	

# Aggregation unit input file

- Aggregation unit is a gene and 20 kb flanking region upstream and downstream of it
- Only subset of 1000K variants used for the workshop were used
- No annotation based filtering was performed on the variants
- Indels are not included

## Header of aggregation unit input file

group_id	chromosome	position	ref	alt
ENSG00000188157.9	1	970546	C	G
ENSG00000242590.1	1	970546	C	G
ENSG00000188157.9	1	985900	C	T
ENSG00000217801.5	1	985900	C	T
ENSG00000242590.1	1	985900	C	T
ENSG00000273443.1	1	985900	C	T

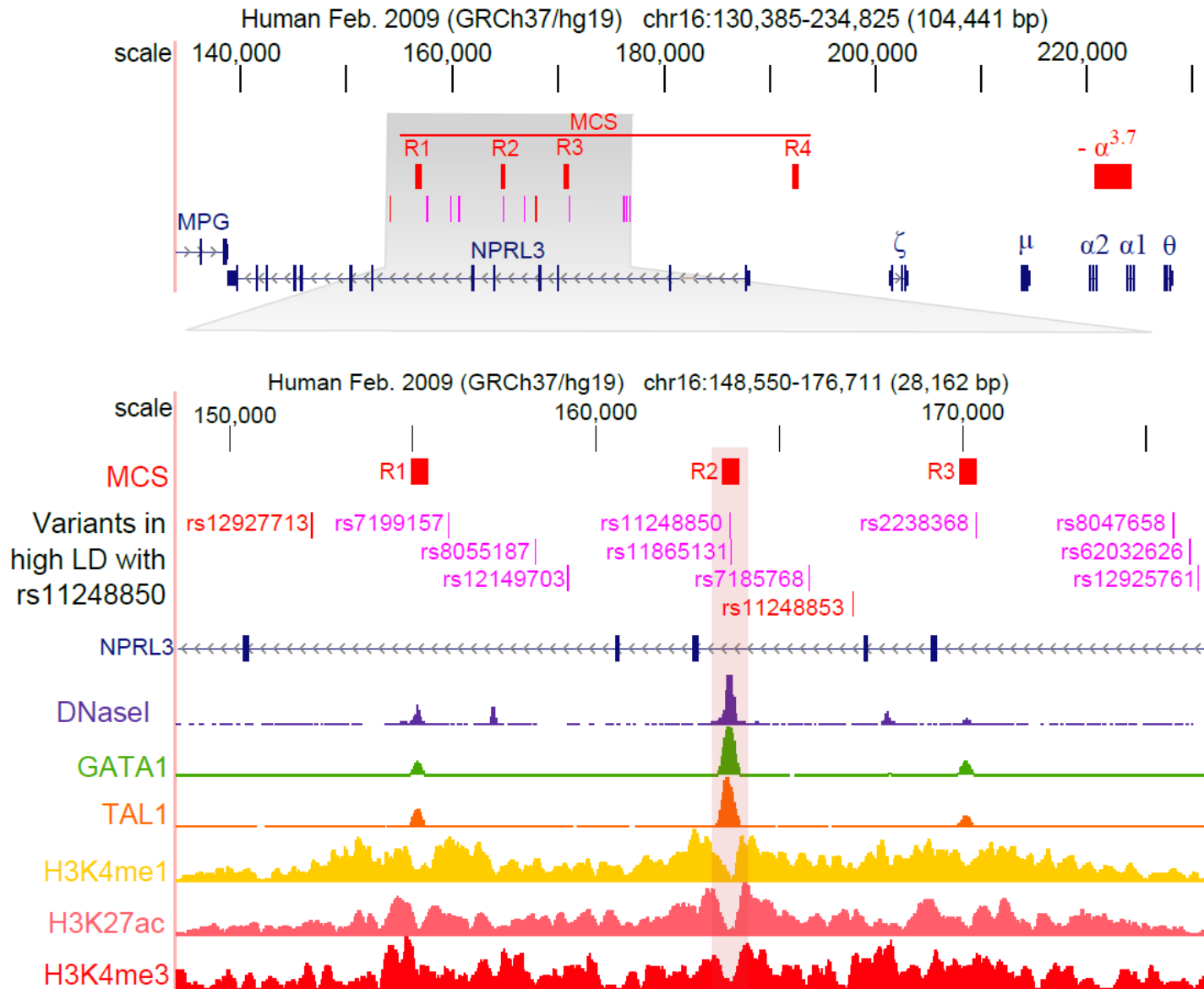
Aggregation unit  
identifier to which  
the variant belongs

Variant information

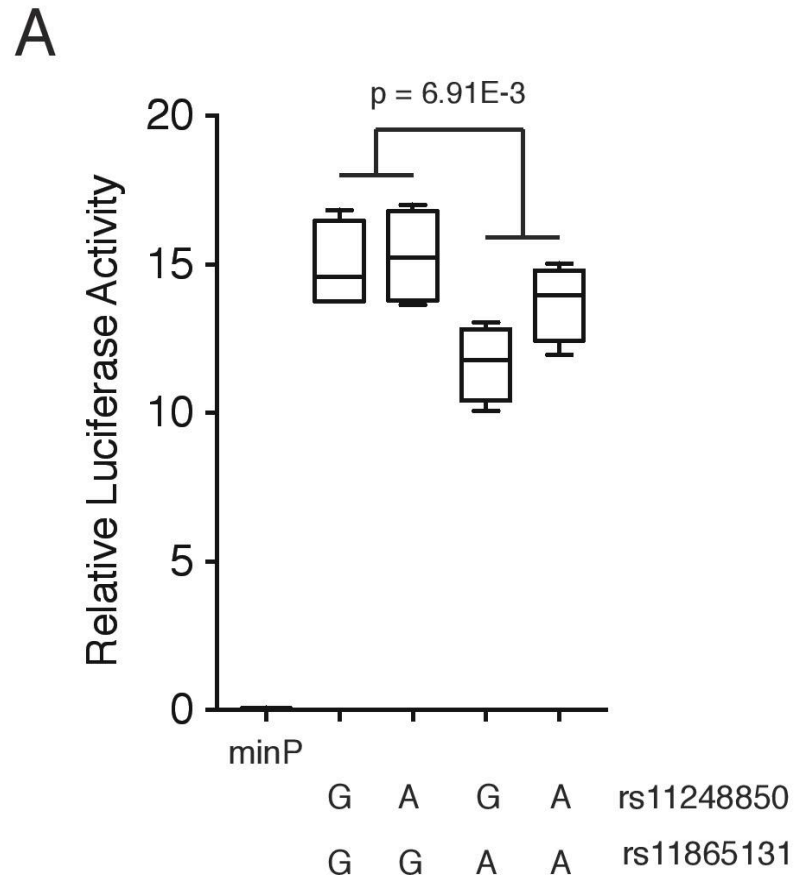
NOTE: A given variant can belong to multiple aggregation unit

# EXTRA SLIDES

# Predicting putative causal variants



# Functional assays confirm allele-specific activity of the predicted causal variants



# Choosing the length of flanking regions for gene-based aggregation units

