

# Figures HW4

Yushi Tang

**1a.** Insert list of genomic regions ([BED](#) format): ([help](#))**1b.** Or [BED](#) file upload:  no file selected**2. Select:**Organism: Assembly: Background: Descriptors: Additional descriptors:  no file selected Compute global pvalue for additional descriptors**Run!** **Reset!**hw4q3\_tumor\_over\_normal\_peak\_74  
55 9.54144

.....

\*\*\*\*\*

9988 regions acquired...

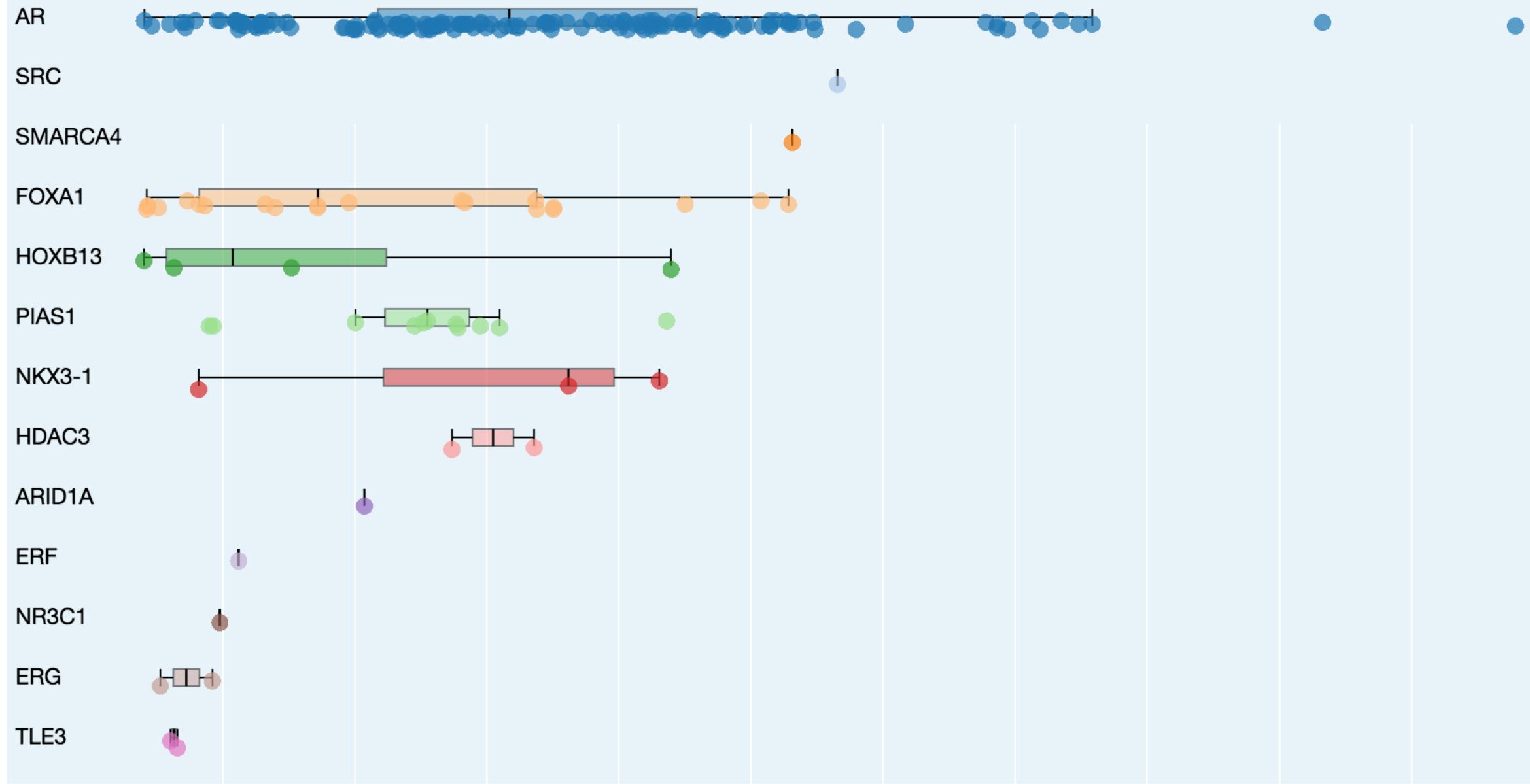
Running Pscan\_Chip...

Please wait... done

**PscAn**  
**chip**Download txt file

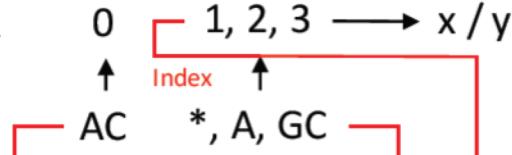
Name	ID	L.PV	L.O/U	G.PV	G.O/U	SP.COR	P.POS	P.POS.PV
<a href="#">Ar</a>	MA0007.3	0		0		0.0191	[-2,8]	5.6E-13
<a href="#">NR3C2</a>	MA0727.1	0		0		0.0119	[-2,8]	3.8E-10
<a href="#">FOXP1</a>	MA0481.2	0		0		0.0126	[1,11]	3.7E-10
<a href="#">FOXF2</a>	MA0030.1	0		0		0.0115	[12,22]	0.0535
<a href="#">Foxo1</a>	MA0480.1	9.3E-274		0		0.002	[-5,5]	4.1E-5
<a href="#">FOXD1</a>	MA0031.1	0		0		-0.0387	[-1,9]	2.6E-6
<a href="#">Foxa2</a>	MA0047.2	0		0		0.0097	[-14,-4]	8.8E-9
<a href="#">FOXA1</a>	MA0148.3	0		0		0.0077	[2,12]	4.5E-9
<a href="#">FOXP2</a>	MA0593.1	0		0		0.0018	[-15,-5]	1.1E-7
<a href="#">FOGX1</a>	MA0613.1	0		0		-0.076	[-16,-6]	3.3E-8
<a href="#">Foxj2</a>	MA0614.1	0		0		-0.024	[-1,9]	2.1E-10
<a href="#">FOXI1</a>	MA0042.2	0		0		-0.0742	[-12,-2]	2.2E-6
<a href="#">FOXL1</a>	MA0033.2	0		0		-0.0739	[3,13]	1.6E-6
<a href="#">FOXO3</a>	MA0157.2	0		0		-0.0142	[-16,-6]	4.6E-6
<a href="#">NR3C1</a>	MA0113.3	0		0		0.0107	[-2,8]	5.2E-11
<a href="#">FOXB1</a>	MA0845.1	0		0		-0.0031	[-15,-5]	0.0263
<a href="#">FOXC1</a>	MA0032.2	0		0		0.0019	[-15,-5]	0.0348
<a href="#">FOXC2</a>	MA0846.1	0		0		0.0046	[-14,-4]	0.0003
<a href="#">FOXD2</a>	MA0847.1	0		0		-0.0746	[-2,8]	1.3E-8
<a href="#">FOXO4</a>	MA0848.1	0		0		-0.0726	[-2,8]	2.0E-8
<a href="#">FOXO6</a>	MA0849.1	0		0		-0.0692	[-2,8]	7.9E-9
<a href="#">FOXP3</a>	MA0850.1	0		0		-0.0723	[-2,8]	1.6E-7
<a href="#">Foxj3</a>	MA0851.1	0		0		0.0027	[3,13]	0.0021

## The most similar samples of your peak sets



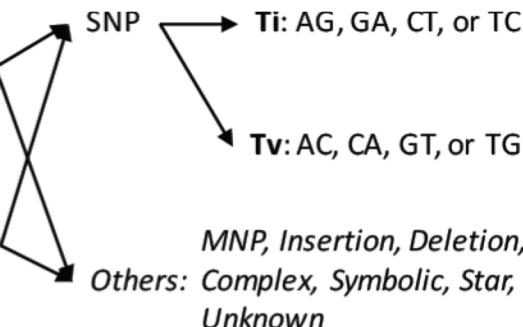
# Feature Calculation

- #Ti and #Tv

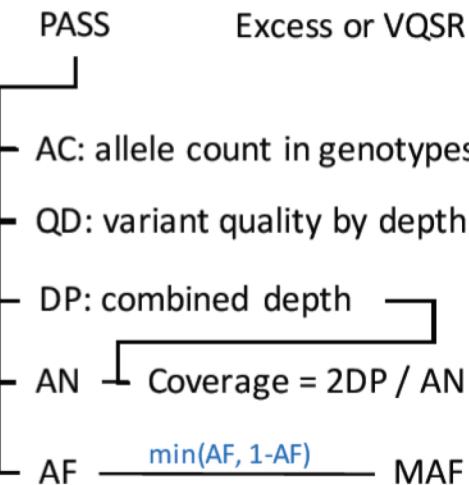


For x:  $x == 0$   
    False → Check the type of mutation from 0 to x  
    True → No transition or transversion happened

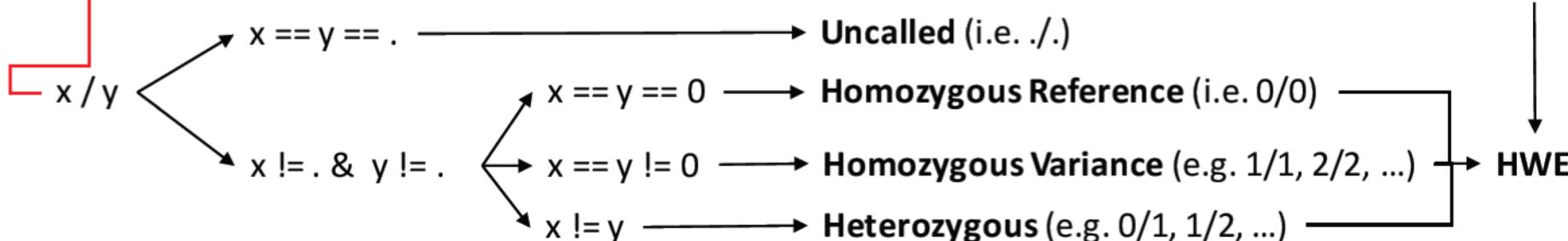
For y:  $y == 0$   
    False → Check the type of mutation from 0 to x  
    True → No transition or transversion happened



## Genotype Quality



- #Het, #HomRef, and #HomVar



# Genotype Quality (GQ & PL)

- Phred-scale genotype likelihoods rounded to the closest integer (PL)

$$PL = -10 \times \log P(\text{Genotype} | \text{Data})$$

- $P(\text{Genotype} | \text{Data})$  is the conditional probability of the Genotype given the sequence Data that we have observed
- Normalization: subtract the value of the lowest PL from all the values, and normalize the values across all genotypes so that the PL value of the most likely genotype is 0
- E.g. Alleles with A (Ref) and T (Alt) with conditional probabilities as:

$$P(AA | \text{Data}) = 0.000001, P(AT | \text{Data}) = 0.000100, P(TT | \text{Data}) = 0.010000$$

Genotype	A/A	A/T	T/T
Raw PL	$-10 \times \log(0.000001) = 60$	$-10 \times \log(0.000100) = 40$	$-10 \times \log(0.010000) = 20$
Normalized PL	$PL(0/0) = 60 - 20 = 40$	$PL(0/1) = 40 - 20 = 20$	$PL(1/1) = 20 - 20 = 0$

- Genotype quality (GQ)
  - The difference between the second lowest PL and the lowest PL



# Genotype Quality (GQ & PL)

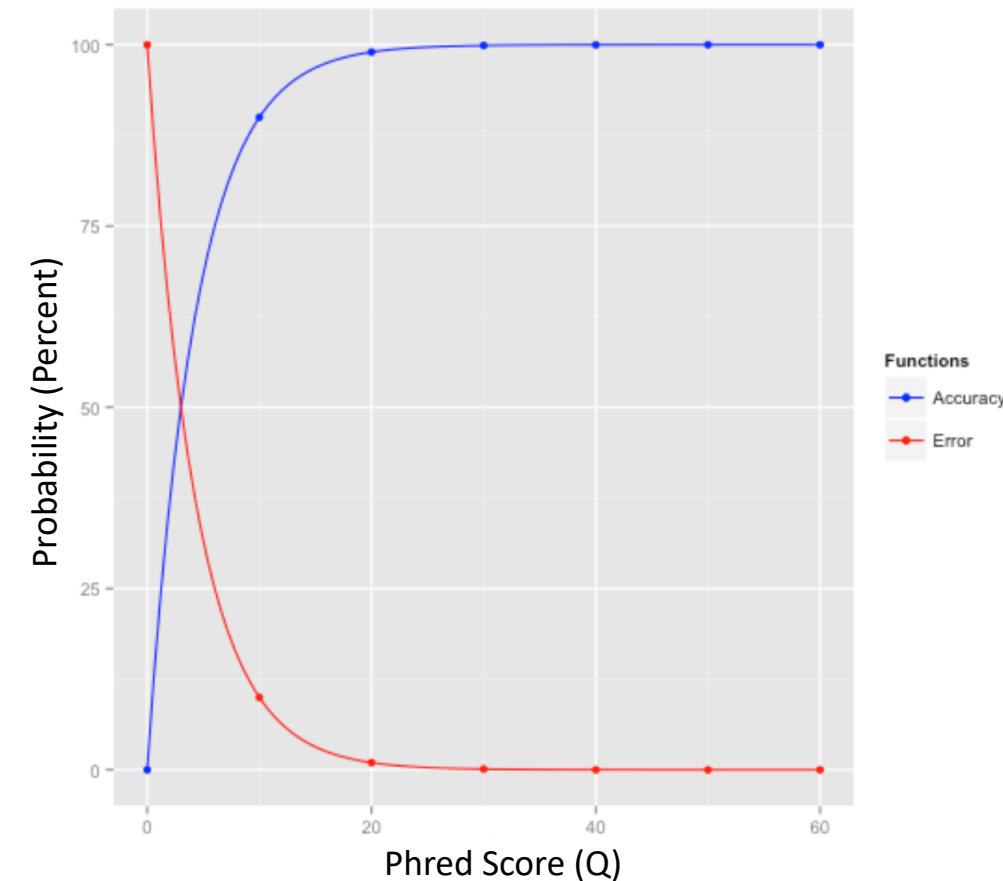
- **Phred-scale in practice: 20 is recommended by GATK as a threshold**
  - The Phred quality score (Q) is logarithmically related to the error probability (E)

$$Q = -10 \times \log E$$

- As for the estimate of accuracy (A), we can take

$$A = 1 - E = 1 - 10^{-\left(\frac{Q}{10}\right)}$$

Phred Quality Score	Error	Accuracy(1-Error)
10	$1/10 = 10\%$	90%
20	$1/100 = 1\%$	99%
30	$1/1000 = 0.1\%$	99.9%
40	$1/10000 = 0.01\%$	99.99%
50	$1/100000 = 0.001\%$	99.999%
60	$1/1000000 = 0.0001\%$	99.9999%



# Data Size Comparison Between Formats

Data Sets	VCF	VCF.BGZ	BCF	GDS
1000 Genome	75 GB	14 GB	12.3 GB	2.6 GB
TOPMed Freeze 5 (n=54,000)	120 TB	22 TB	987 GB	143 GB
GSP CCDG Freeze 1 raw (n = 22,000)	436 TB	55 TB	3.6 TB	450 GB
GSP CCDG Freeze 1 after QC (n = 22,000)	21 TB	2 TB	126 GB	21 GB
Compression	1	0.09	0.006	0.001

The files in different formats contain the same data