# Association Studies

## Normal (Wild-Type) Haemoglobin

**Haemoglobin**
Normal (globular)

**Red Blood Cell**
Round (biconcave)

**Blood Vessels**
Free flowing

RBC

## 'Sickle Cell' Haemoglobin

**Haemoglobin**
Clumped (fibrous)

**Red Blood Cell**
(sickle-shaped)

**Blood Vessels**
Forms clots / blockages

Sickle Cell

[https://www.23andme.com/ancestry-composition-guide/](https://www.23andme.com/ancestry-composition-guide/)

## Ancestry Composition:

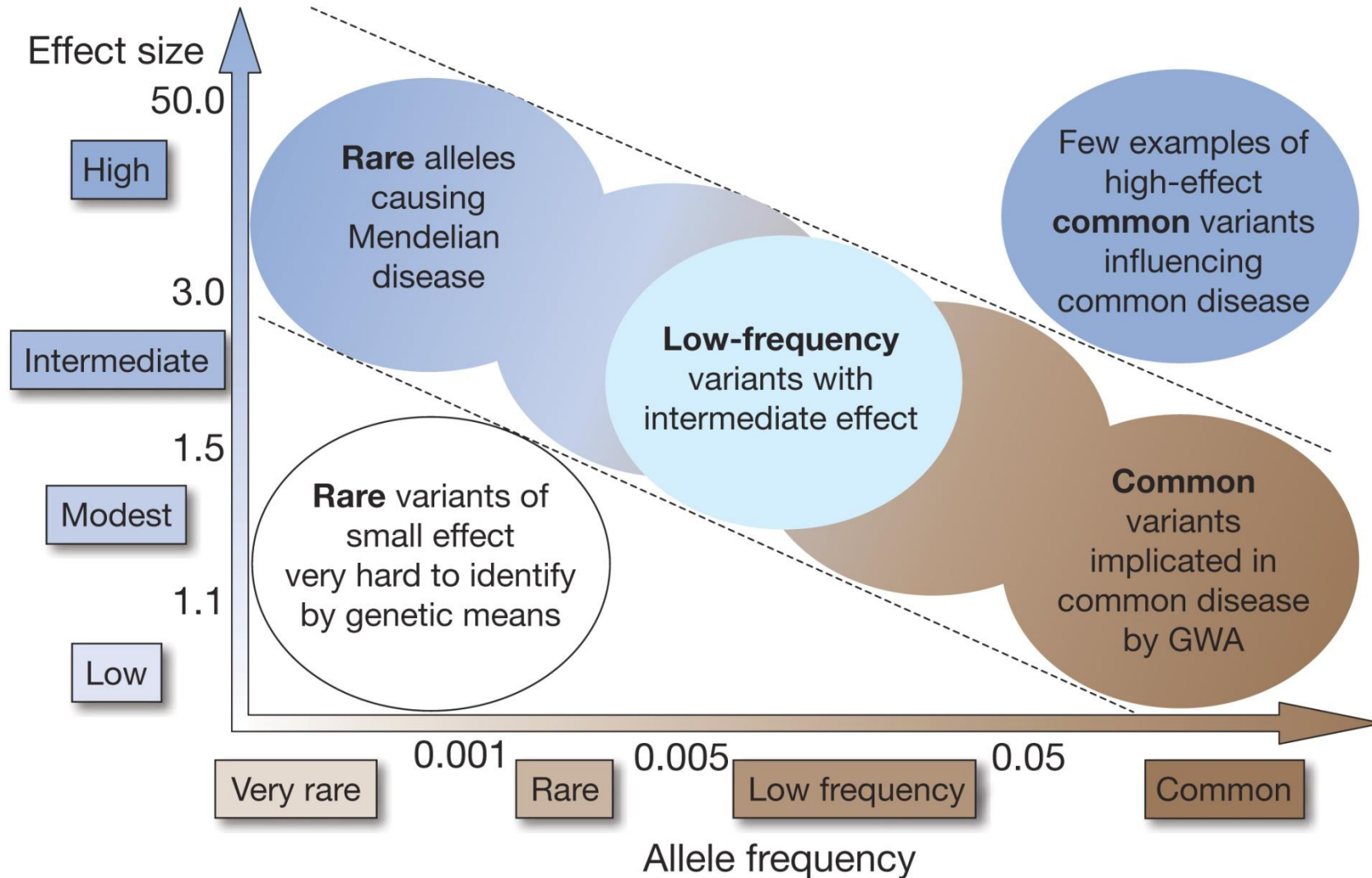### 23andMe's State-of-the-Art Geographic Ancestry Analysis

23andMe's Ancestry Composition report is a powerful and well-tested system for analyzing ancestry based on DNA, and we believe it sets a standard for rigor in the genetic ancestry industry. We wrote this document to explain how our analysis works and to present some quality-control test results. Note: This document goes into specifics for the current version of Ancestry Composition, offered to customers on the V5 platform. For customers on previous platforms, click here.

Your Ancestry Composition report shows the percentage of your DNA that comes from 45 populations. We calculate your Ancestry Composition by comparing your genome to those of over 10,000 people with known ancestry. When a segment of your DNA closely matches the DNA from one of the 45 populations, we assign that ancestry to the corresponding segment of your DNA. We calculate the ancestry for individual segments of your genome separately, then add them together to compute your overall ancestry composition.

# Learning objectives

- Describe the differences and the pros and cons of sequencing vs genotyping.

- Calculate and interpret odds ratios in case/control genetic association studies.

- Interpret quantitative trait association studies.

- Understand role for imputation.

# Genetic Variation and Disease



Manolio et al. Nature 2009; 461: 747-753.

# Genetic data collection

- TaqMan Polymerase chain reaction (PCR)
  - Targeted, low throughput.
  - Detect deletions and structural variations.
- Genotyping chip
  - Targeted locations, high throughput.
  - Detects single, *a priori* locations.
- Sequencing
  - Collects all bases, increasingly high throughput.
  - Identify novel variants.
  - Analyzing data more intensive

# TaqMan PCR to identify variants

# Genotyping technologies (low-throughput)



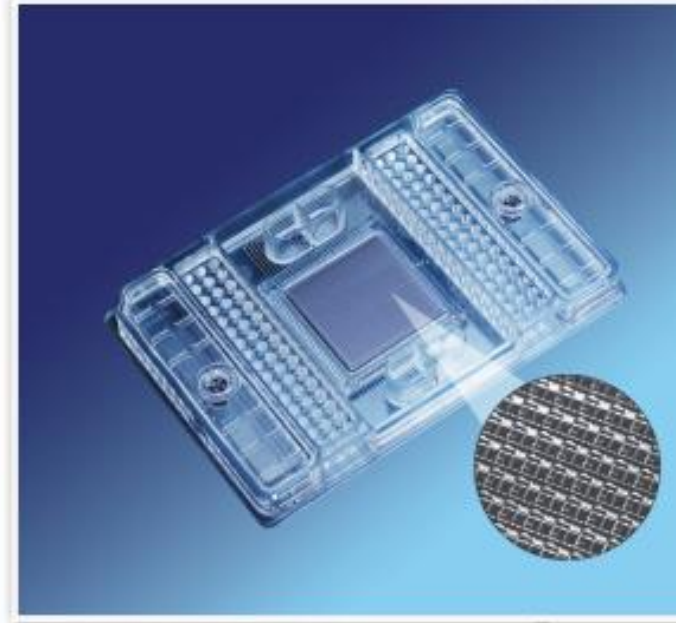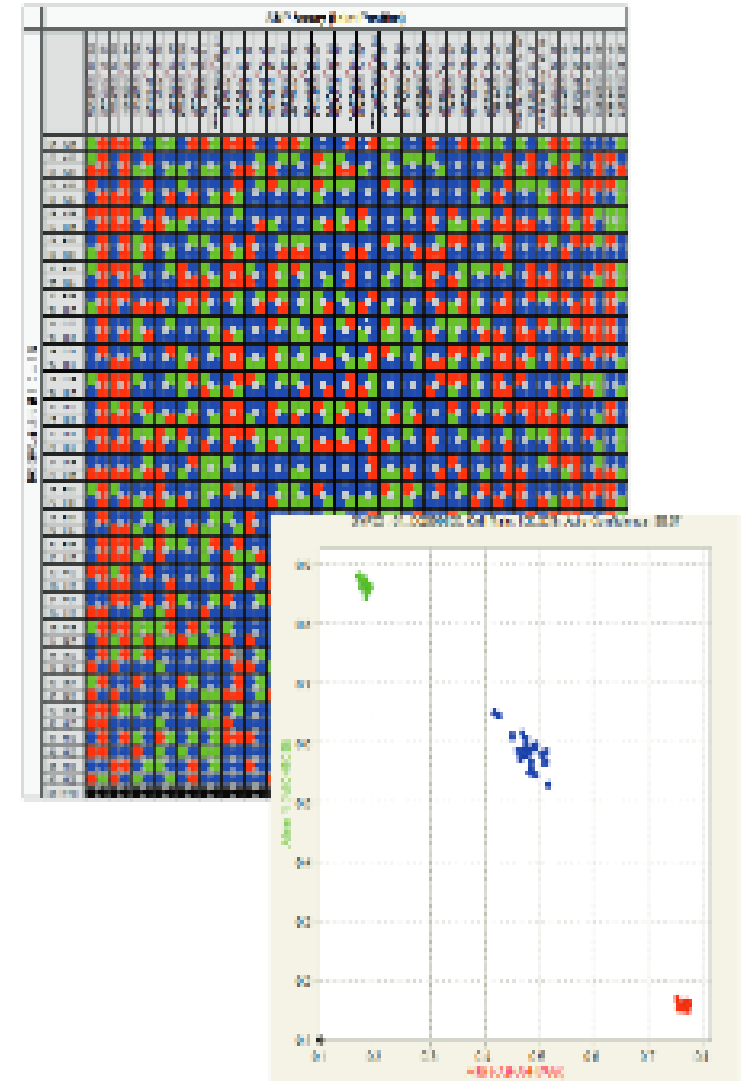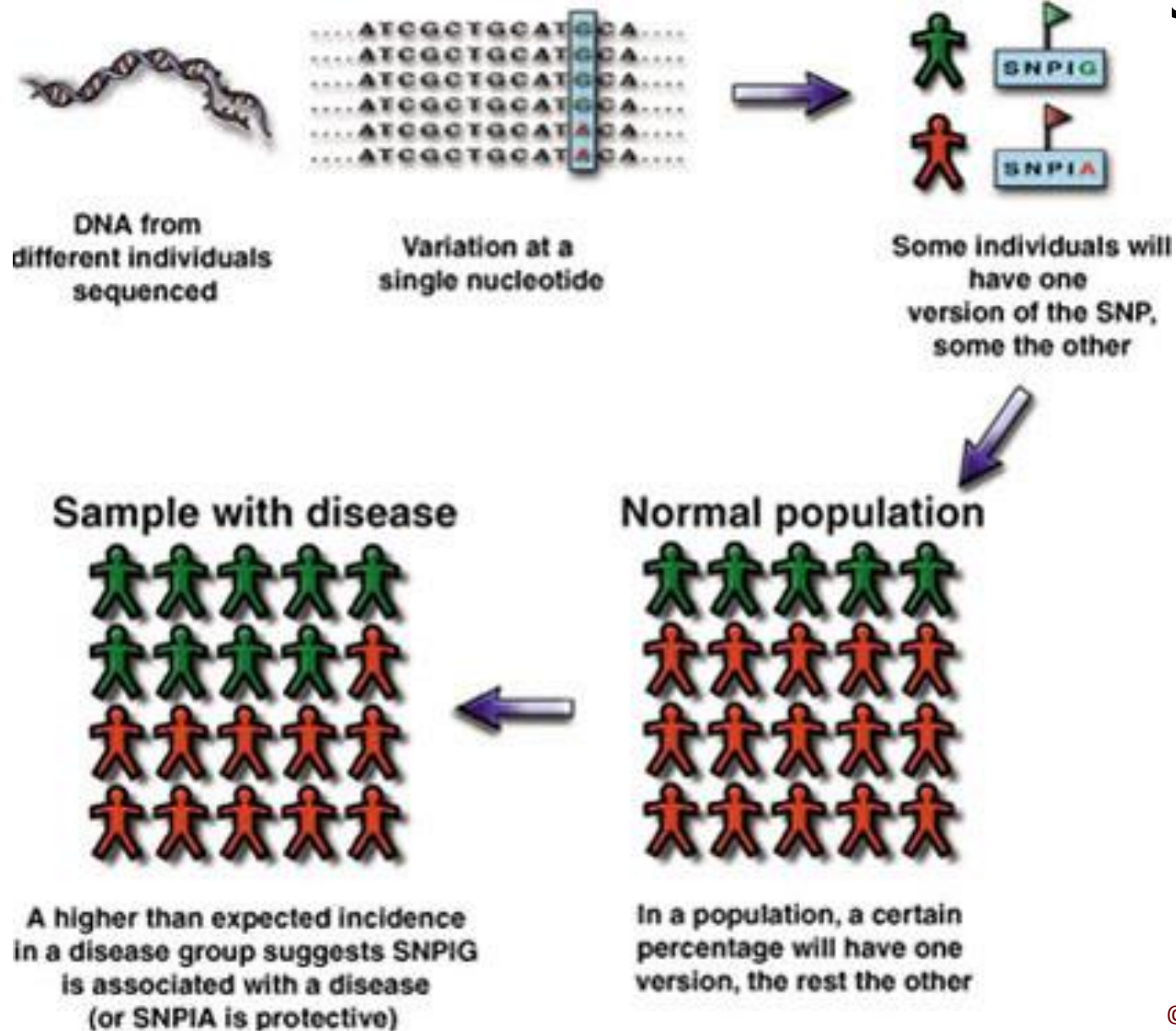| Illumina | SNPlex | Sequenom | TaqMan |
|----------|--------|----------|--------|
| 1500 - 300 SNPs | 400 - 40 SNPs | 40 - 5 SNPs | 10 - 1 SNPs |

# Chip Genotyping

Why we like SNPs:

- Abundant in the genome
- Easy to measure



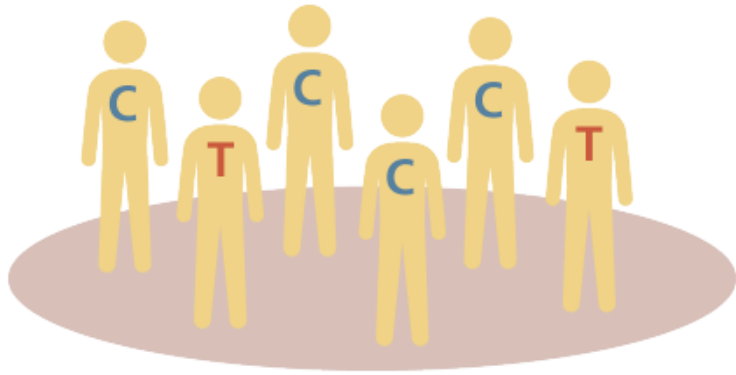Microfluidics, 96 samples x 96 assays, DNA probes with fluorescent markers.

Fluidigm platform

# Genetic association studies using SNPs



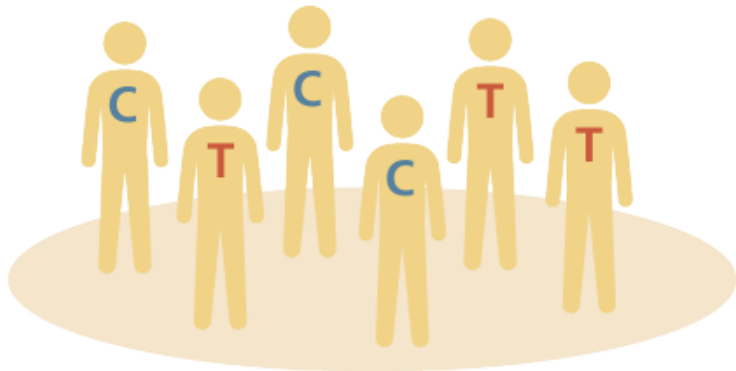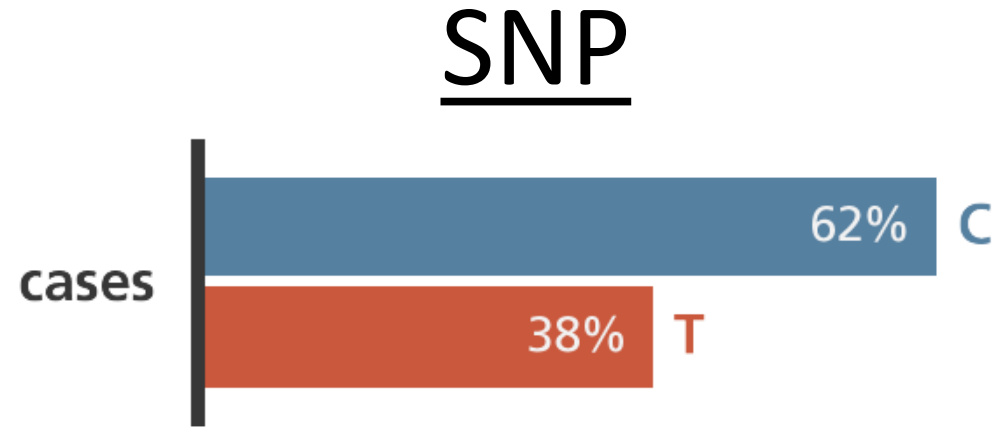© Gibson & Muse, A Primer of Genome Science

# Association studies

- Determine if a particular genetic feature (exposure) co-occurs with a trait (disease) more often than would be expected by chance.

- Binary: Calculate 'odds' of an outcome occurring.
  - Framed as an 'odds ratio', the odds of an outcome after an exposure (genotype) in relation to the odds of an outcome without the exposure (reference genotype).

- Continuous: calculate change in an outcome for every unit increase of an exposure.

cases (n=1,000)
people with heart disease

controls (n=1,000)
people without heart disease

## SNP

| | | |
|---|---|---|
| cases | 62% C | |
| | 38% T | |

| | | |
|---|---|---|
| controls | 49% C | |
| | 51% T | |

# Odds ratio

The odds ratio is our measure of association for a case-control study. It tells us whether and how much an exposure increases the likelihood of our outcome of interest. We need to look at two things:

**The estimate** -- the odds ratio itself. How big in the connection between an exposure and an outcome? Are those with an exposure more likely to have the outcome?

**The p-value** -- how certain are we that the odds ratio didn't just happen by chance?

# Association testing in case-control studies

| | | Disease status | | |
|---|---|---|---|---|
| | | Cases | Controls | Total |
| **Genotype** | M | a | b | a+b |
| | m | c | d | c+d |
| **Total** | | a+c | b+d | |

measure of events out of all possible events (Ratio) vs ratio of events to non-events (Odds)

$$RR = \frac{\text{Risk of event in the Treatment group}}{\text{Risk of event in the Control group}} = \frac{a/(a+b)}{c/(c+d)}$$

$$OR = \frac{\text{Odds of event in Treatment group}}{\text{Odds of event in Control group}} = \frac{a/b}{c/d} \quad :$$

If an outcome occurs 10 out of 100 times, the risk is 10%
But the odds is 10/90 = 11.1%

# Association testing in case-control studies

| | | Disease status | | |
|---|---|---|---|---|
| | | Cases | Controls | Total |
| **Genotype** | M | a | b | a+b |
| | m | c | d | c+d |
| **Total** | | a+c | b+d | |

1) Calculate the odds of the disease with the genotype and without the genotype

Odds that the M genotype occurs in a case: $= \dfrac{a}{b}$

Odds that the m genotype occurs in a case: $= \dfrac{c}{d}$

# Association testing in case-control studies

| | | Disease status | | |
|---|---|---|---|---|
| | | Cases | Controls | Total |
| **Genotype** | M | a | b | a+b |
| | m | c | d | c+d |
| **Total** | | a+c | b+d | |

2) Calculate Odds Ratio (OR) as the odds that genotype M occurs in a case divided by the odds that genotype m occurs in a case.

$$\frac{a/b}{c/d} = \frac{ad}{bc}$$

OR $= \dfrac{ad}{bc}$

# Association testing in case-control studies

| | | Disease status | | |
|---|---|---|---|---|
| | | Cases | Controls | Total |
| **Genotype** | M | a | b | a+b |
| | m | c | d | c+d |
| **Total** | | a+c | b+d | |

Odds that the M allele occurs in a case $= \dfrac{a}{b}$

Odds that the m allele occurs in a case $= \dfrac{c}{d}$

The Odds Ratio (OR) is the odds that M occurs
in a case divided by the odds that m occurs in a case:

$$OR = \dfrac{ad}{bc}$$

**$H_0$: OR = 1   (no association)**

**OR > 1   indicates increased odds**

**OR < 1   indicates decreased odds
(protective)**

# Confidence intervals for odds ratios

| | | Disease status | |
|---|---|---|---|
| | | Cases | Controls |
| **Genotype** | M | a | b |
| | m | c | d |

$$\text{OR} = \frac{a/b}{c/d} = \frac{ad}{bc}$$

$$\text{s.e}(\log(\text{OR})) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

Confidence interval: $e^{\log(OR) \pm z_{\alpha/2} \times s.e(\log(OR))}$

Lower limit of 95% confidence interval: $e^{\log(OR) - 1.96 \times s.e}$

Upper limit of 95% confidence interval: $e^{\log(OR) + 1.96 \times s.e}$

# Calculate– odds ratio and 95% confidence interval

|  | Cases | Controls | Total |
|---|---|---|---|
| TT+TC | 158 | 392 | 550 |
| CC | 20 | 86 | 106 |
| Total | 178 | 478 | 1656 |

$$OR = \frac{ad}{bc}$$

$$s.e(\log(OR)) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

# Odds ratio calculations – odds ratio itself

|  | Cases | Controls | Total |
|---|---|---|---|
| TT+TC | 158 | 392 | 550 |
| CC | 20 | 86 | 106 |
| Total | 178 | 478 | 1656 |

$$OR = \frac{158 \times 86}{392 \times 20} = 1.7332$$

$$s.e.(log(OR)) = \sqrt{\frac{1}{158} + \frac{1}{392} + \frac{1}{20} + \frac{1}{86}}$$

# Odds ratio calculations – confidence intervals

|  | Cases | Controls | Total |
|---|---|---|---|
| TT+TC | 158 | 392 | 550 |
| CC | 20 | 86 | 106 |
| Total | 178 | 478 | 1656 |

$$OR = \frac{158 \times 86}{392 \times 20} = 1.7332$$

$$s.e.(log(OR)) = \sqrt{\frac{1}{158} + \frac{1}{392} + \frac{1}{20} + \frac{1}{86}}$$

lower limit 95% confidence interval:

$$= exp(log(OR) - 1.96 \times s.e.(log(OR)))$$

$$= exp(log(1.7332) - 1.96 \times 0.2665) = 1.03$$

Upper limit 95% confidence interval: 2.92

# Odds ratio calculations – odds ratio itself

|  | Cases | Controls | Total |
|---|---|---|---|
| TT+TC | 158 | 392 | 550 |
| CC | 20 | 86 | 106 |
| Total | 178 | 478 | 1656 |

OR = 1.7

Turn this result into a sentence about effect of T allele in thyroid cancer.

# Odds ratio calculations – odds ratio itself

| | Cases | Controls | Total |
|---|---|---|---|
| TT+TC | 158 | 392 | 550 |
| CC | 20 | 86 | 106 |
| Total | 178 | 478 | 1656 |

OR = 1.7

Turn this result into a sentence about effect of T allele in thyroid cancer.

The odds of developing thyroid cancer are 1.7x times greater with an T allele compared to without an T allele.

# Why do we even use odds and odds ratios???

The odds ratio allows us to calculate the associations between an exposure and an outcome without needing the frequency of the exposure in the general population

>(very useful to rare exposures, such as rare diseases).
>(we'd have to sample A LOT of people to get a true population picture and even pick up one or two cases of the disease)

The log(odds) allows us to transform this weird variable into a linear form, which is easier for us to fit to models and interpret the output.

# Why do we use Log odds 5:26 - 8:42

# Often use logistic regression for case-control analyses

Allows you to adjust for relevant factors
- Population stratification, age, sex, matching variables etc

$$\ln\left(\frac{p}{1-p}\right) = \alpha + \beta_1 \mathbf{g} + \beta_2 x_1 + \ldots + \beta_{k+1} x_k \quad \text{(g is genotype, } x_1, \ldots x_k \text{ are covariates)}$$

Coefficients are estimated using maximum likelihood estimation (MLE)

- $\ln\left(\frac{p}{1-p}\right)$ = log odds of an outcome
- Test $H_0$: $\beta_1 = 0$ (likelihood ratio test, wald test, score test)
- The odds ratio is OR=$e^{\beta_1}$
- $\beta_1$ = SNP effect (log(OR)) ➔ e$^{\beta_1}$ = OR

# Logistic regression output

```
> Association<- glm(binaryPhenotype~HLA.B5701,family=binomial(link="logit"),data=AbacavirData)
> summary(Association)

Call:
glm(formula = binaryPhenotype ~ HLA.B5701, family = binomial(link = "logit"),
    data = AbacavirData)

Deviance Residuals:
    Min       1Q    Median       3Q       Max
-1.3770  -1.3770    0.3349   0.9902    2.4478

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)   -2.944      1.026  -2.870  0.00410 **
HLA.B5701P     3.402      1.051   3.236  0.00121 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 138.63  on 99  degrees of freedom
Residual deviance: 114.76  on 98  degrees of freedom
AIC: 118.76

Number of Fisher Scoring iterations: 5
```

**c=log(odds of allergy)**

**a**

**x**

**b**

# Common models of penetrance

Effect

Effect

Effect

**Recessive**
Genotype coding: 0,0,1

**Dominant**
Genotype coding: 0,1,1

**Additive**
Genotype coding: 0,1,2

AA    AC    CC

AA    AC    CC

AA    AC    CC

Effect = mean of continuous trait or log(OR) of binary trait

# Continuous outcome genetic association

- Linear regression (instead of logistic)
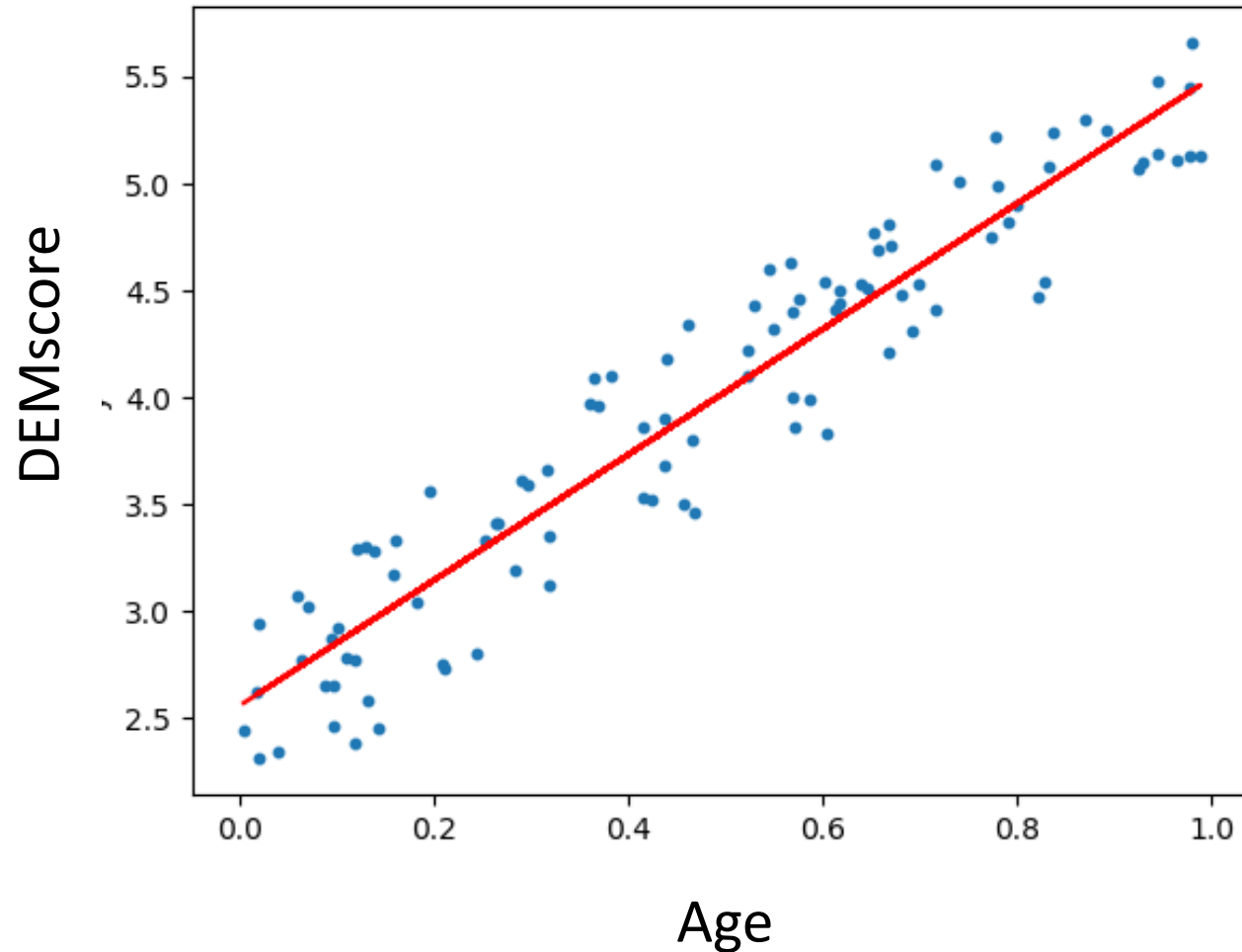- Additive coding of SNP (0,1,2) most common

$$Y = \alpha + \beta * SNP + X$$

- $\beta$ = SNP effect (for every SNP, unit increase in outcome)
- SNP = covariate coded (0,1,2)
- X = additional covariates (e.g. sex, study, age, PCs from population stratification)

# Multivariate analyses

# One predictor, one outcome

# Multivariate analysis

# Importance of setting your reference allele

Odds ratio when AA is reference: $\frac{2}{3} / \frac{1}{3} = \frac{2}{3} * 3 = 2$

**The odds of the outcome are 2x more likely among those with CC genotype compared to among those with the AA genotype.**

Odds ratio when CC is reference. $\frac{1}{3} / \frac{2}{3} = \frac{1}{3} * \frac{3}{2} = 0.5$

**The odds of the outcome are ½ as likely among those with AA genotype compared to among those with the CC genotype.**

**These are the saying the same thing! But the language matters.**

# Always know and be purposeful on your reference

In epidemiology, the reference group always matters.

**Exposure** (gene allele reference)

**Outcome** (some outcomes have no "direction") brown vs black hair

**Population** (other factors are always involved, i.e. age, diet, access to care).

# Surrogate endpoints



1. Type 2 diabetes → HbA1c → Macrovascular complications

2. Osteoporosis → (Bisphosphonates) → Bone mineral density → Fracture risk

# Picking an endpoint

- Surrogate endpoints: an endpoint that in itself means nothing, but gives information about an important endpoint.
    - More proximal in the biological pathway.
        - Time to detect and/or intervene.
        - Can detect earlier and collect more people.
    - Easier to measure, especially if an outcome often results in death.
    - Monitor progress and change in risk.
    - Cheaper to measure and conduct study.
- Problems with surrogate endpoints
    - Misclassification -- loss of precision

We could also have turned Bone mineral density into a binary outcome based on whether the measure was below the threshold for high fracture risk:



Is the BMD more of less than 500g/cm?

# Quantitative vs categorical outcomes

## Quantitative

- Does not rely on subjective labels.
- Often more likely to detect differences.
- Interpretation: increase in unit change of phenotype per unit change in risk factor.

## Binary

- Must decide cutpoint.
- More straightforward message for action.
- Interpretation: increase in odds of phenotype per unit change in risk factor.

# Interpretable cutpoints -> aids policy development

Genotyping platforms can vary by studies, how can we combine data or get more genotyping data than we start with?

# We can use LD in our studies: tagSNPs

# We can use LD in our studies: Imputation



HapMap or 1,000 Genomes

Reference haplotypes

```
0  0    1   1 1 0  0    1    1    0 0    0   1 1    1
0  0    0   0 0 1  1    1    0    1 1    1   0 0    1
1  1    1   1 1 0  0    0    1    0 0    0   0 0    0
1  0    1   1 0 0  0    1    1    1 1    1   0 0    1
```

Cases and controls typed on SNP chip

```
1  ?    ?   ? 2 ?  0    ?    ?    ? ?    0   1 ?    1
1  ?    ?   ? 1 ?  0    ?    ?    ? ?    ?   0 ?    0
0  ?    ?   ? 1 ?  1    ?    ?    ? ?    1   0 ?    1
1  ?    ?   ? 2 ?  0    ?    ?    ? ?    0   1 ?    1
?  ?    ?   ? 2 ?  0    ?    ?    ? ?    0   0 ?    0
1  ?    ?   ? 1 ?  1    ?    ?    ? ?    1   0 ?    ?
0  ?    ?   ? 2 ?  0    ?    ?    ? ?    0   1 ?    1
1  ?    ?   ? 1 ?  1    ?    ?    ? ?    1   1 ?    2
```

Study genotypes

Hirschhorn & Daly. *Nature Reviews Genetics 2005,*
http://mathgen.stats.ox.ac.uk/impute/impute_v2.html

# Imputation

Due to LD, we can compare haplotypes between a "reference" panel and our study and thereby guess genotypes

**Study Individual:** T A G G T **?** T G C C T A **?** C G T

**Reference Panel Individual:** T A G G T **A** T G C C T A **G** C G T

https://mathgen.stats.ox.ac.uk/impute/impute_v2.html

Genotyping

Person 1    ---T------G---A
Person 2    ---T------G---A
Person 3    ---T------C---A
Person 4    ---A------G---T
Person 5    ---T------C---A
Person 6    ---A------G---T

Match genotypes to a reference

GGCTATTTTGGGAA
CGCTATATACCCAT
GGCAATTTAGCGAT
GCCTATATACGGAA

Can you impute the missing bases?

# Imputation

- Cost efficient
  - Can assess more SNPs than we genotyped (tagSNPs)

- Allows us to keep our sample size
  - Fill in missings for already genotyped SNPs

- Allows us to combine data from existing platforms and different studies that genotype different SNPs

# Imputation

- We can infer genotypes for SNPs we didn't genotype (or failed in the lab)
  - **Input:** 550,000 SNPs in 10,000 individuals

  - **Reference panel:** 2,504 individuals from the 1000 Genomes project (>80M markers)

  - **Output:** Imputed data for >80M markers for your 10,000 individuals
    - In practice, we exclude markers that were only seen once in 1000Genomes so we end up with ~47M markers)

# Assessing SNPs across genotyping platforms

|  | HumanHap | Affy 6.0 | OmniExpress |
|---|---|---|---|
| **HumanHap** | 459,999 | 126,959 | 260,661 |
| **Affy 6.0** |  | 668,283 | 168,223 |
| **OmniExpress** |  |  | 565,810 |

\* 75,285 markers are on all 3 platforms

Lindström, PLoS One 2017

# Imputation for studying SNPs across platforms

Illumina SNPs

Affymetrix SNPs

Overlap SNPs

# Imputation for studying SNPs across platforms

1000G SNPs

Illumina SNPs

Affymetrix SNPs

Overlap SNPs

GWAS genotypes

Reference panel 1 (e.g., study-specific)

Reference panel 0 (e.g., 1000 Genomes)

—— Known genotypes

- - - Genotypes being imputed

① Impute panel 0-specific variants into panel 1.

② Impute panel 1-specific variants into panel 0.

③ Use merged reference panel to impute untyped variants in GWAS dataset.

# Imputation

- The imputation quality score $r^2$ measures how well a SNP was imputed.
  - Ranges between 0 and 1.
  - A quality score of $r^2$ on a sample of $N$ individuals indicates that the amount of data at the imputed SNP is approximately equivalent to a set of perfectly observed genotype data in a sample size of $r^2N$.
  - Typically, a cut-off of 0.30 or so will flag most of the poorly imputed SNPs, but only a small number (<1%) of well imputed SNPs. Caveat: This is not true for rare SNPs

# Imputation

- Factors that affect imputation quality:

    - Number of genotyped SNPs in your data

    - Size of reference panel

    - Similarity in genetic ancestry between reference and study samples

    - Allele frequency

What if we don't know what variants to test or they are too rare to impute?

# Sequencing vs Genotyping: Discovery

Array imaged to analyze signal at all beads

**Genotyping**:

- Common variants (>5% allele frequency)
- large cohorts (cheaper)
- to identify regions of the genome associated with an outcome
- less computationally demanding to get a person's alleles.

**Sequencing**:

- Rare variants
- Discover new variants in individuals or small samples (compare children and parents)
- very detailed data
- to add variants across the same gene in studying an effect.

# Genotyping Output



Li, Nat Comm 2014

# Genotype cluster plot for rare variants



Auer, Nat Genet 2014

# Sequencing to identify rare variants.

Same variant shared by individuals in a small group.

Multiple variants in the same gene in individuals with the same condition.

Variants unique to an individual in an important gene.

# Sequencing Technologies

- Sanger sequencing uses real time PCR
  - 99.99% accuracy
  - Used for high-accuracy reads of smaller regions
- Next Gen sequencing sequences many segments at once
  - Also called: massively parallel sequencing
  - High throughput
  - Used for multi-gene reads and larger samples

# Sequencing output

# Sequencing alignment and depth

Depth: The number of times one basepair is sequenced

# Nanopore sequencing

Works on a single DNA strand. No PCR amplification; No chemical labeling.

Feeds DNA through a *very* tiny hole, sends an electric current, and determines the DNA base based on how the current flows.



Nanopore sequencing of DNA

# Sequence Assembly

CTCGCGCGAT

ACCCTCG

GCGATAG

ACTTAATAC

ACCCTCGCGC

GCGATAGACTTA

# Sequence Assembly

ACCCTCGCGCGATAGACTTAATAC

---

CTCGCGCGAT

ACCCTCG

GCGATAG

ACTTAATAC

ACCCTCGCGC

GCGATAGACTTA

# Sequence Assembly: Other Considerations

- Assembly type: De novo or mapping sequence assembly

- Read length: usually 100-700bp

- Read depth: 30 is gold standard

# Sequence Assembly: Why Read Depth Matters

ACCCTCGCGCGATAGACTTAATAC

ACCCTCG
ACCCGCGCGC
CTCGCGCGAT
GCGATAGACTTA
GATAG
ACTTAATAC

| Name | Lengt | Av. q | Quality graph | Bases |
|---|---|---|---|---|
| SRR000702.28 | 36 | 39 | | CACATAGGAGTCCAGAACACTGCTGCTGAGGTATAA |
| SRR000702.28 | 36 | 34 | | TGCCTGCCTGAGGACTCTGGTGCTGGAGGCTGTCTT |
| SRR000702.28 | 36 | 39 | | CCTTGGCCTCTCAAAACGCTGAGATTACAGGCGTGA |
| SRR000702.28 | 36 | 29 | | CACATATACACACCTCCACATACACACAGATCGG |
| SRR000702.28 | 36 | 37 | | CATGGGCCTGTAGGATTAGATAAGCATACTTGCTAT |
| SRR000702.28 | 36 | 34 | | CACTGGGGCTTTCATCGGACGCTGTGTCTCACCGCG |
| SRR000702.28 | 36 | 33 | | CAGCACTGAGTTTCTGAGAGAGTGGCCAGCTGGGCT |
| SRR000702.28 | 36 | 30 | | TTGTATTTGGCAAGGGGTTGCTTGTTATAGCTTGTT |
| SRR000702.28 | 36 | 38 | | CAGGAGAAGGGAAATGTGGGTTGGAAGCTTTAATTG |
| SRR000702.28 | 36 | 27 | | CATATAAAACCCTCTTCCCCTTTCAACACACTTAAT |
| SRR000702.28 | 36 | 39 | | CTCGGCTCACTGCAAACTCTGCTTCCCAGGTTCATG |
| SRR000702.28 | 36 | 1 | | TNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN |
| SRR000702.28 | 36 | 31 | | CAGTGTTGCTTGCTTCTGTTTTACATGTACTAGTAG |
| SRR000702.28 | 36 | 40 | | CAATACACATCTACCGACACACACACTCATACACAC |
| SRR000702.28 | 36 | 37 | | CAAGAGGCATGGGGGATGTGCTCTATCCTGTTTTGT |
| SRR000702.28 | 36 | 36 | | CATTCCATTCTATTGCATTCCATTCTATTCTGTTTA |
| SRR000702.28 | 36 | 40 | | CCATTCCATTCCATTCCATTCCATTCCATTCCATTC |
| SRR000702.28 | 36 | 29 | | ATCCATTTACATAGCAAATGGCTGGATGTGCCCTTC |
| SRR000702.28 | 36 | 6 | | GGCAGGAGCTCCCCATGTGCTGCAACAGCTTCCTAA |
| SRR000702.28 | 36 | 38 | | CCCACGGTGTCCATAAGTGGAGTCAATGCCTCTGAA |
| SRR000702.28 | 36 | 28 | | TATATCACACACACATTTATACACTCAAACTGTTT |
| SRR000702.28 | 36 | 38 | | CACTTGATTTTTGAGCCTTATAATAAGGCTAGAGAG |
| SRR000702.28 | 36 | 34 | | TCTCCCCACAGATGAGCAGCAGCTGCTCAGGGCTGA |
| SRR000702.28 | 36 | 39 | | CATGCACCGCAACATTCAGCTAGTATTTTTATTTTT |
| SRR000702.28 | 36 | 38 | | CAAACTATTCACACACAAACTCTACACACATATAAA |
| SRR000702.28 | 36 | 37 | | TGCATAATCTTGGCTCACTGCAACCTCCACTTCCAG |

# Phred Q scores: probability of incorrect call

p = 10^(-Q/10)

Sequence ID    Sequenced Read

```
@SRR081708.237049/1
GGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
+
I>:<9>>>:=:>>?<>:@?>;==@@@>?=AAA<>=A@?6>4B=<>>.@>?<@;?#############
```
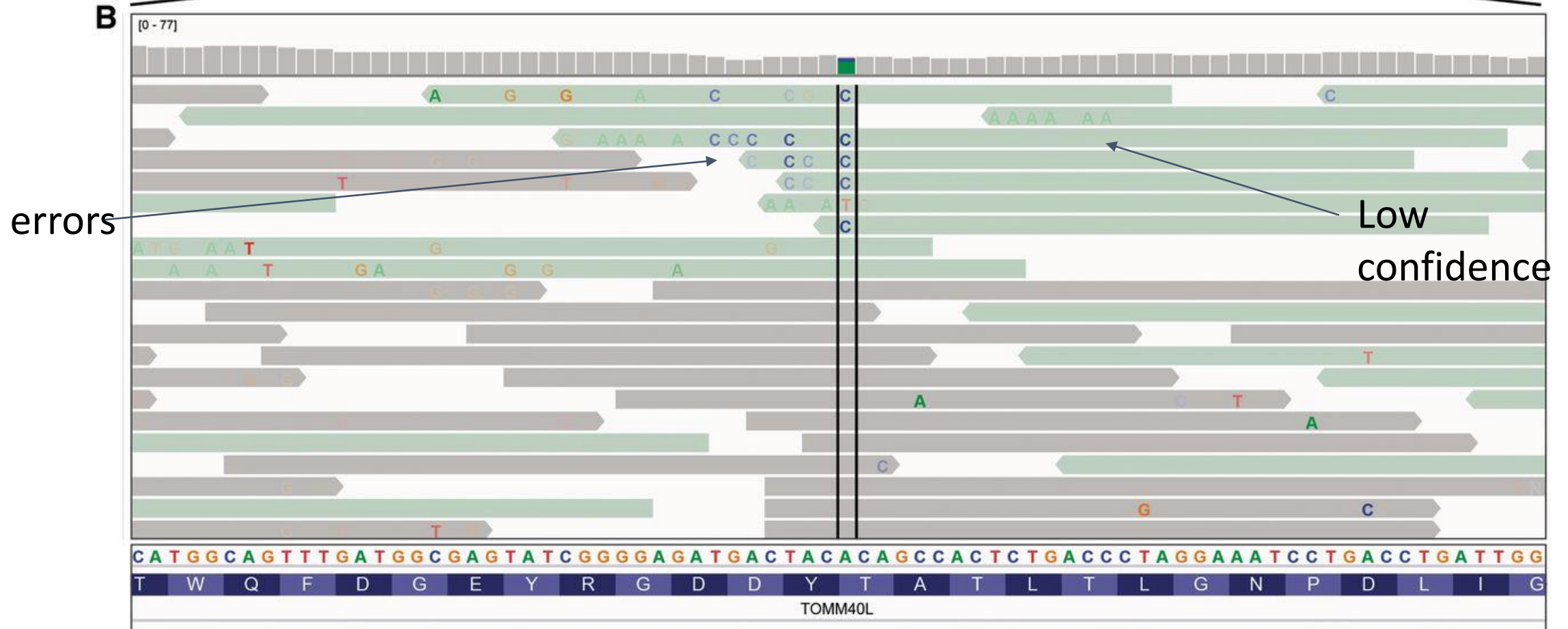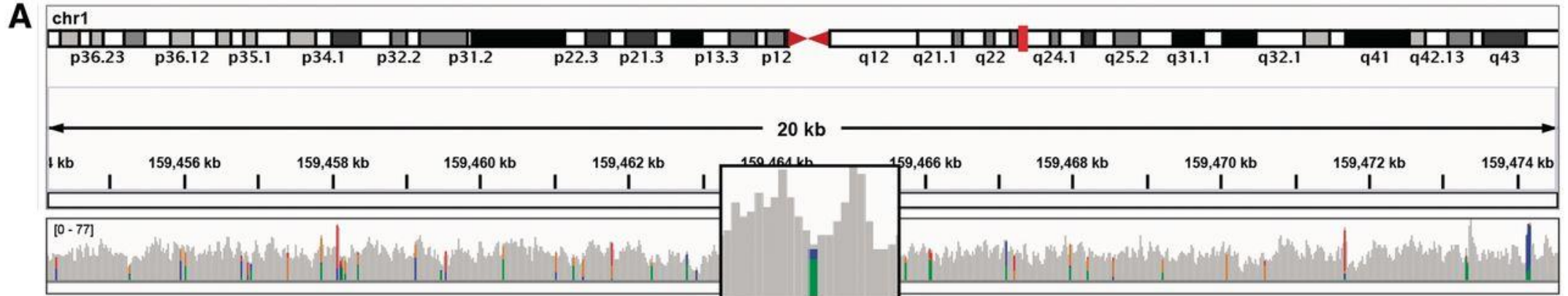
Blank    Quality Score

**Phred quality scores are logarithmically linked to error probabilities**

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |
| 60 | 1 in 1,000,000 | 99.9999% |

Higher the number, the better

errors

Low confidence

# You have a new variant. Now what??

If they change the amino acid structure, algorithms to predict pathogenicity:

**SIFT: Sorting Intolerant From Tolerant.** Based on sequence homology across species and chemical properties of amino acids. If it is the same, it must be important. Closer to zero = worse.

**Polyphen-2:** Based on protein structure and function predictions. Such as where in the protein the change occurs. Closer to 1 = worse.

Can have very different answers! Prediction is hard.

## Multiple sequence alignment

Shown are 75 amino acids surrounding the mutation position (marked with a black box). An interactive version of the complete alignment is also available.

## 3D Visualization
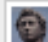
PDB/DSSP Snapshot 03-Jan-2012 (78304 Structures)



**EntryID:** 2IF1

**ChainID:** A

**Residue:** Leu72

**Identity:** 100.0%

**Overlap:** 100.0% (113 aa)

Zoom into mutation    Reset view    View size: + −

Jmol script terminated

Web design & development:

# Variant Effect Predictor ❓

**Species:**

👤 Human (Homo sapiens) ✏️

Assembly: GRCh37.p13 *(If you are looking for VEP for Human GRCh38, please go to GRCh38 website🔗.)*

**Name for this job (optional):**

**Input data:**

**Either paste data:**

**Examples:** Ensembl default, VCF, Variant identifiers, HGVS notations, SPDI

**Or upload file:** [ Choose File ] No file chosen

**Or provide file URL:**

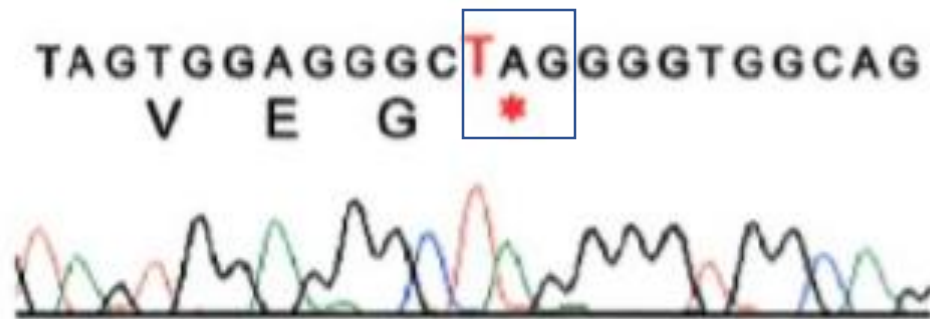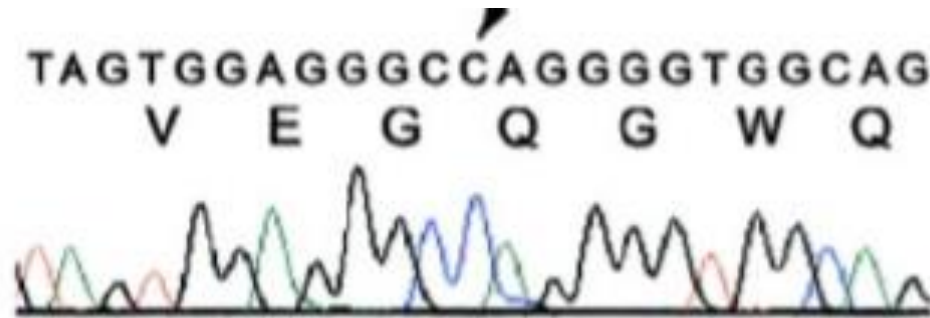**Transcript database to use:**

- ⦿ Ensembl/GENCODE transcripts
- ○ Ensembl/GENCODE basic transcripts
- ○ RefSeq transcripts

# Sequencing in Founder Populations

- Osteoporosis is a disease in elderly people, resulting in decreased bone density
  - Many treatments may be carcinogenic
- A treatment for osteoporosis discovered by identifying SOST gene implicated in sclerosteosis in Dutch Afrikaner population
  - Autosomal recessive disorder

All but one patient with sclerosteosis in the 22-family sample shared
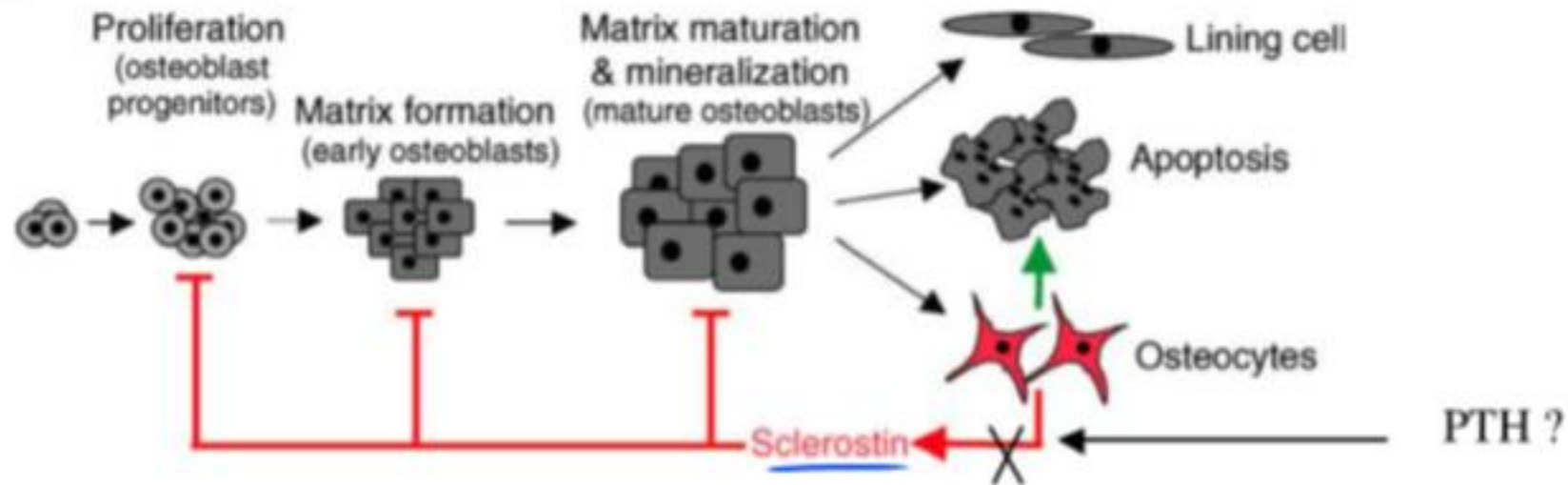the same SNP



Results in a premature stop codon

# -> Drug Development

- A drug was developed to deplete or inhibit sclerostin in people with osteoporosis

# Summary

- Genetic data can be collected through genotyping or sequencing.

- Odds ratios give the odds of an outcome in relation to a reference.

- Linear and logistic regression allow adjustment for other factors.

- Imputation leverages linkage disequilibrium to estimate data not collected.