

# Genetics and Bioinformatics

**Kristel Van Steen, PhD<sup>2</sup>**

**Montefiore Institute - Systems and Modeling**

**GIGA - Bioinformatics**

**ULg**

**[kristel.vansteen@ulg.ac.be](mailto:kristel.vansteen@ulg.ac.be)**

## **Complicating factors in bioinformatics**

### **1 Trait heterogeneity in GWAs**

**Single traits association tests**

### **2 Confounding**

**3.a Epidemiology**

**3.b GWAs (population structure)**

### **3 Multiple testing**

**Locus heterogeneity**

### **4 Multiple studies**

**Meta-analysis**

## **5 When variants become rare – sparse data**

**Customizing GWAs for rare variants association analyses – from GWAs to Sequence Analyses**

## **6 When effects become non-independent**

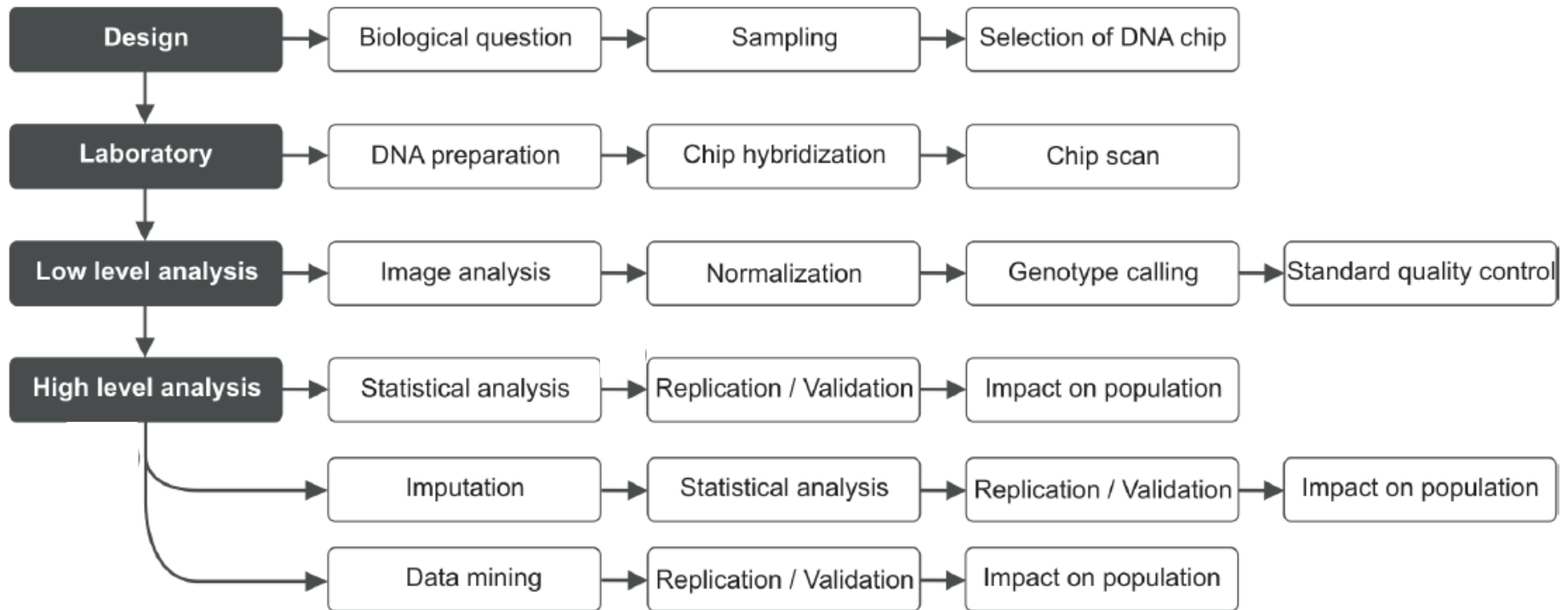
**Biological vs statistical epistasis (future class)**

## 5 When variants become rare – sparse data – expanding GWAs



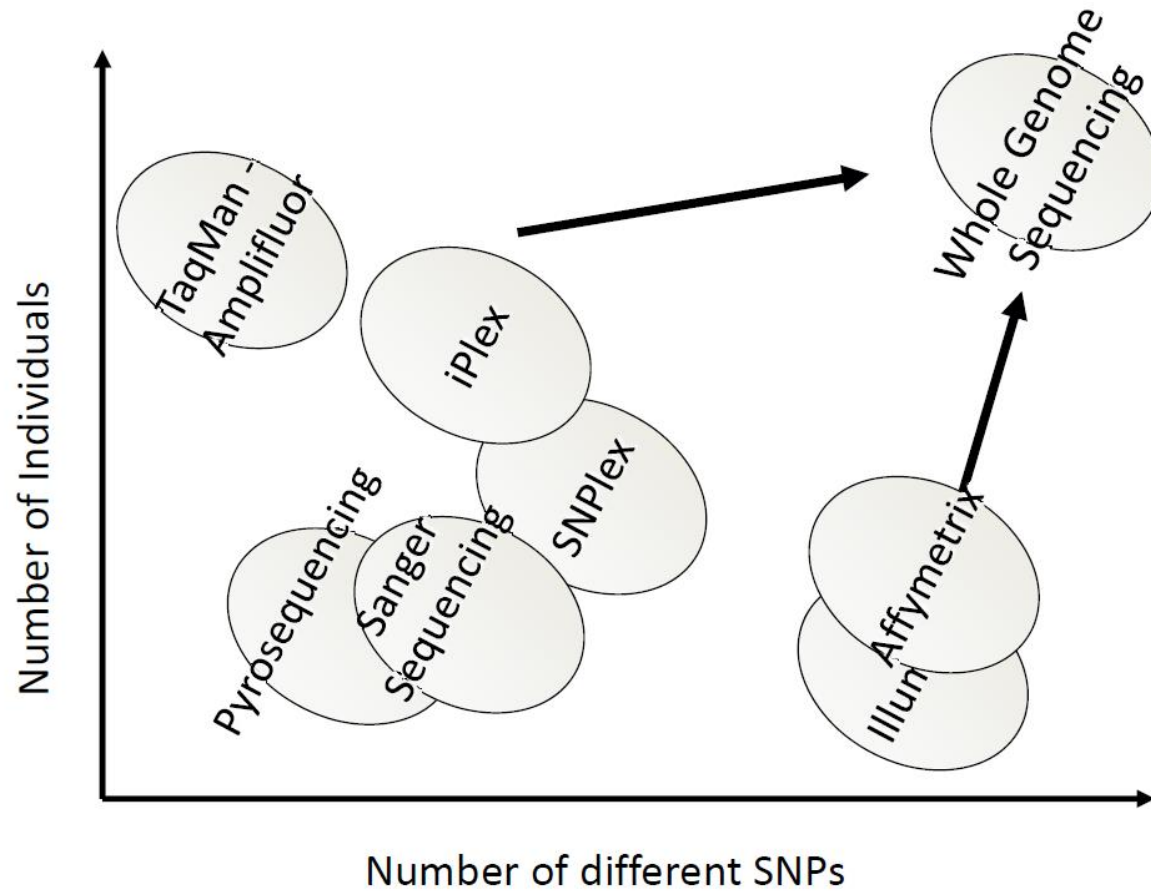
(slide Doug Brutlag 2010)

## Detailed flow of a genome-wide association study

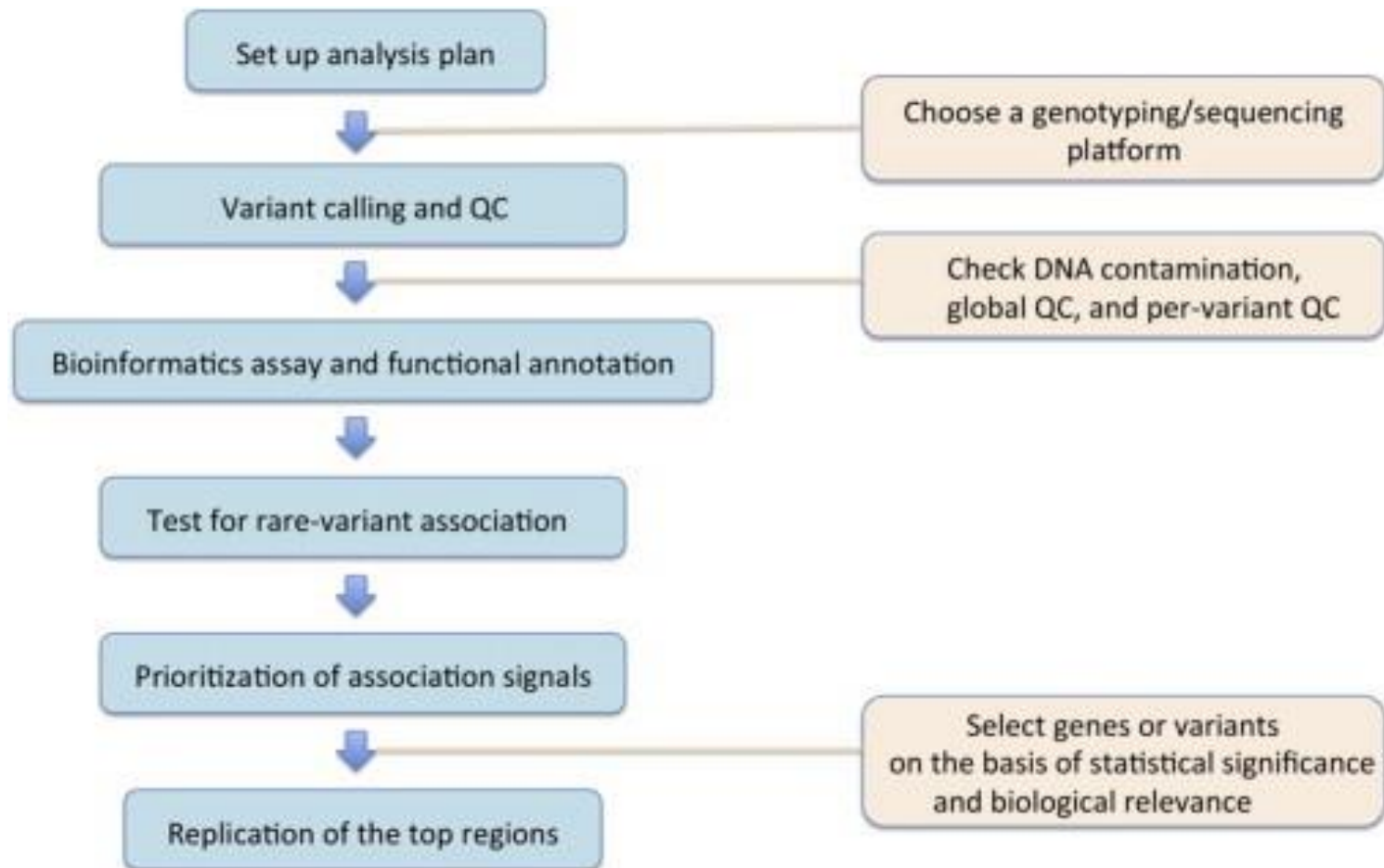


(Ziegler 2009)

## From arrays to sequence data

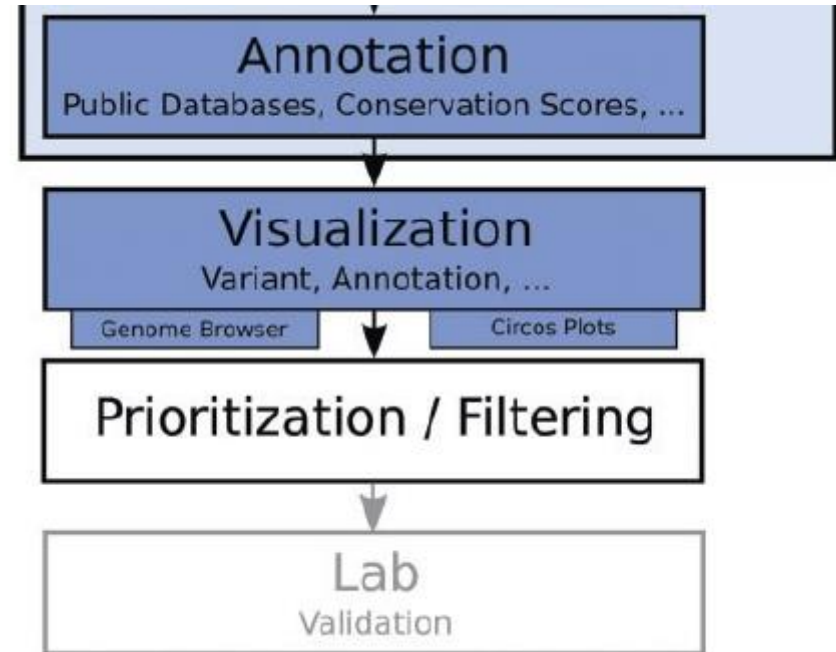
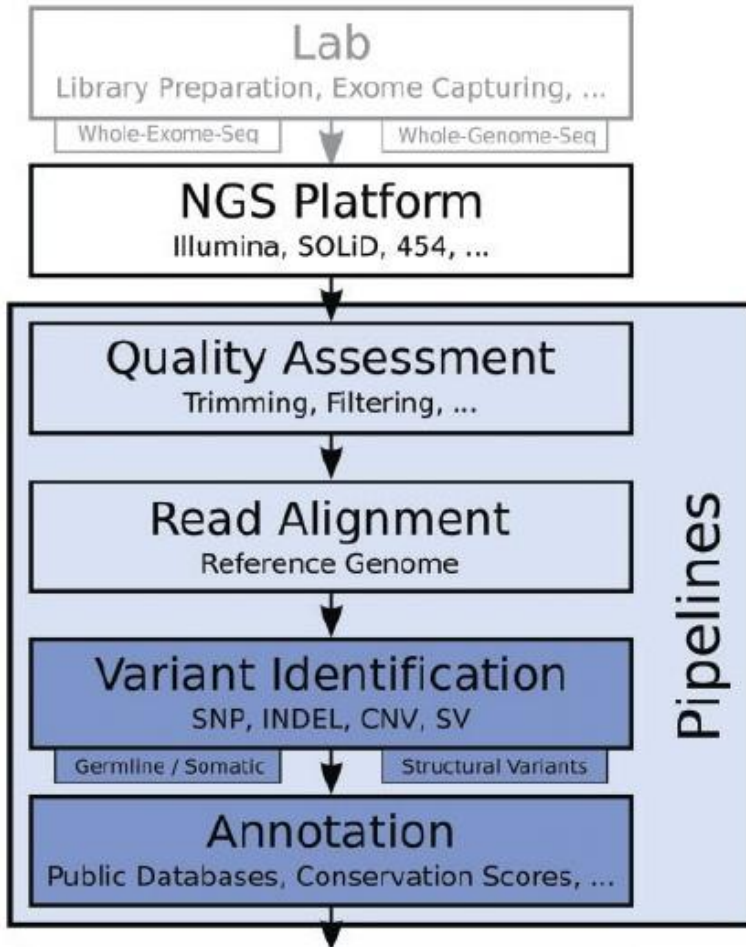


## Work flow genome-wide association study with sequence data



(Lee et al. 2014)

# Common workflow for whole-exome and whole genome sequencing



(Pabinger et al. 2013)



# **A primer on rare variant association testing**

Fan Li  
BIOSTAT 790

January 28, 2016

## Outline

**Overall goal:** to understand the necessity of identifying rare variants and the proposed statistical methods for rare-variant association testing.

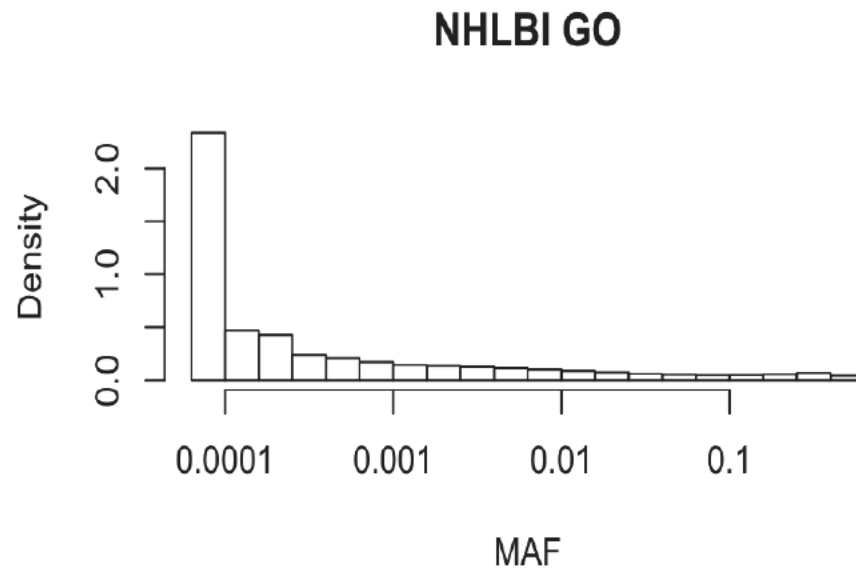
- Rationale for studying rare variants (complex trait)
- Sequencing and study designs
- Rare-variant association tests
- Summary

## Common vs rare variants

- MAF: frequency at which the least common allele occurs in population
- Common variants:  $MAF \geq 5\%$
- Low frequency variants:  $0.5\% \leq MAF < 5\%$
- Rare variants:  $MAF \leq 0.5\%$

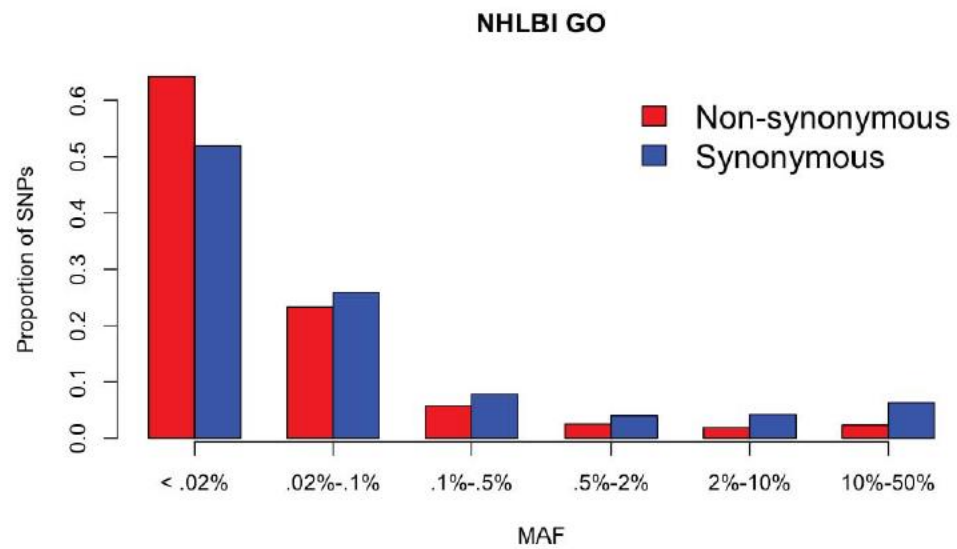
## Why rare variants?

- Most of human variants are rare



## Why rare variants?

- Functional variants tend to be rare



## Why rare variants?

- Further, since common variants explained limited variation in the trait . . .
- Some argued rare variants could explain additional trait variability
- Advancement of sequencing technology (NGS), reduction in cost

## Challenges

- Require cost-effective study designs to genotype many individuals

It can be shown that at least  $\log(1 - \theta) / [2 \log(1 - MAF)]$  individuals are needed to observe a variant with no less than  $\theta$  chance. For  $\theta = 99.9\%$ , we have

MAF(%)	10	1	0.1	0.01
Minimum sample size	33	344	3453	34537

- Classical single-variant tests, developed for common variants detection, are underpowered
- Multiple testing

## Challenges

- A variant – genetic association test implies filling in the table below and performing a chi-squared test for independence between rows and columns

	<b>AA</b>	<b>Aa</b>	<b>aa</b>
<b>Cases</b>			
<b>Controls</b>			

Sum of entries = cases+controls
------------------------------------

- How many observations do you expect to have two copies of a rare allele?  
Example: MAF for a = 0.001 → expected aa frequency is 0.001 x 0.001 or 1 out of 1 million



- **In a chi-squared test of independence setting** (comparing two variables in a contingency table to see if they are related):

When  $MAF \lll 0.05$  then some cells above will be sparse and large-sample statistics (classic chi-squared tests of independence) will no longer be valid. This is the case when there are less than 5 observations in a cell

$$X^2 = \sum_{all\ cells\ i} \frac{(O_i - E_i)^2}{E_i} \quad (\text{contrasting Observed minus Expected})$$

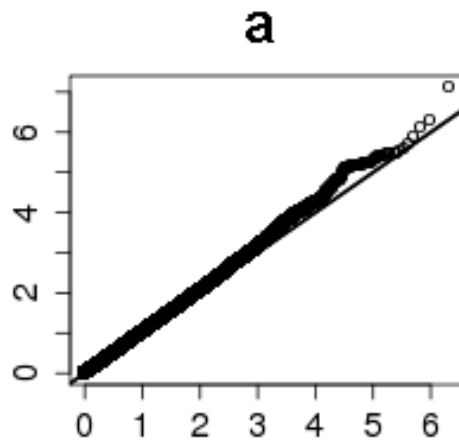
- **In a regression framework:**

The minimum number of observations per independent variable should be 10, using a guideline provided by Hosmer and Lemeshow (Applied Logistic Regression, one of the main resources for Logistic Regression)

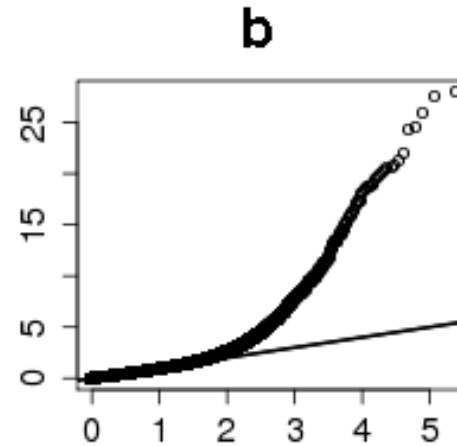
# Increased false positive rates

Q-Q plots from GWAS data, unpublished

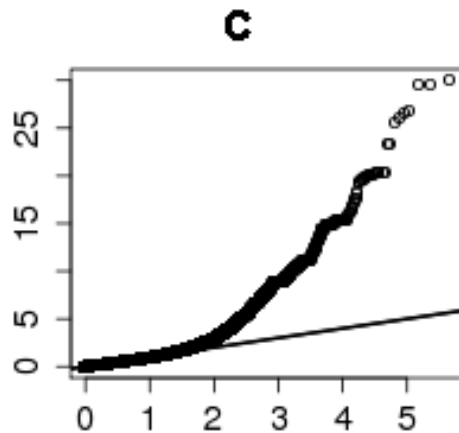
N= $\sim$ 2500  
MAF $>$ 0.03



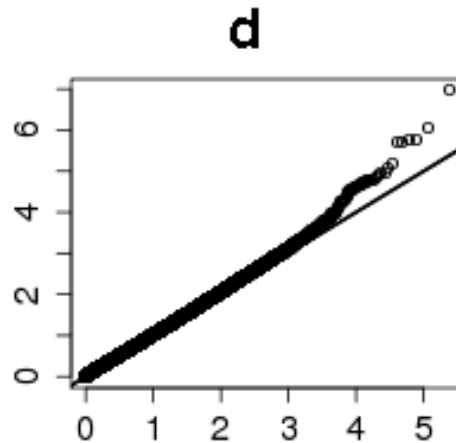
N= $\sim$ 2500  
MAF $<$ 0.03



N= $\sim$ 2500  
MAF $<$ 0.03  
Permuted



N=50000  
MAF $<$ 0.03  
Bootstrapped



## **Remediation: do not look at a single variant at a time, but collapse**

- Rationale for aggregation tests
  - Alpha level of 0.05, corrected by number of bp in the genome=  $1.6 \times 10^{-11}$
  - One needs VERY LARGE samples sizes in order to be able to reach that level, even if you find “the variant”.
- Remedy = aggregate / pool variants
  - Requires specification of a so-called “region of interest” (ROI)
  - A ROI can be anything really:
    - Gene
    - Locus
    - Intra-genic area
    - Functional set

## Remediation: design alternatives to deep sequencing

- Low-depth whole-genome sequencing: sequencing depth refers to the average number of reads that cover each base; **limited accuracy**
- Exome sequencing: **limited to exome**
- High-priority region sequencing: **limited to the target region**
- ...

In summary, either the sequenced range or accuracy is compromised

## How do aggregation tests for (rare) multiple variants work?

- Region-based: gene, regulatory region
- Identify multiple genetic variants within a region
- Evaluate the joint effects of these variants while adjusting for covariates
- **Caution:** These tests rely on assumptions for genetic model (e.g.: mode of inheritance), and the power depends on the true **disease model**  $h(\mu(Y))$

## Corresponding regression models

- $n$  subjects ( $i = 1, \dots, n$ )
- $m$  variants in a region
- Allele counts in a region  $\mathbf{G}_i = (G_{i1}, \dots, G_{im})'$ , ( $G_{ij} = 0, 1, 2$ )
- $q$ -dimensional covariates  $\mathbf{X}_i$  (age, gender, PC scores etc.)
- The disease model is given by a GLM

$$h(\mu(Y_i)) = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}_i + \boldsymbol{\beta}'\mathbf{G}_i \quad (1)$$

- Now the interest is in the null of **no genetic-region effect**:

$$H_0 : \boldsymbol{\beta} = (\beta_1, \dots, \beta_m)' = \mathbf{0}_{m \times 1}$$

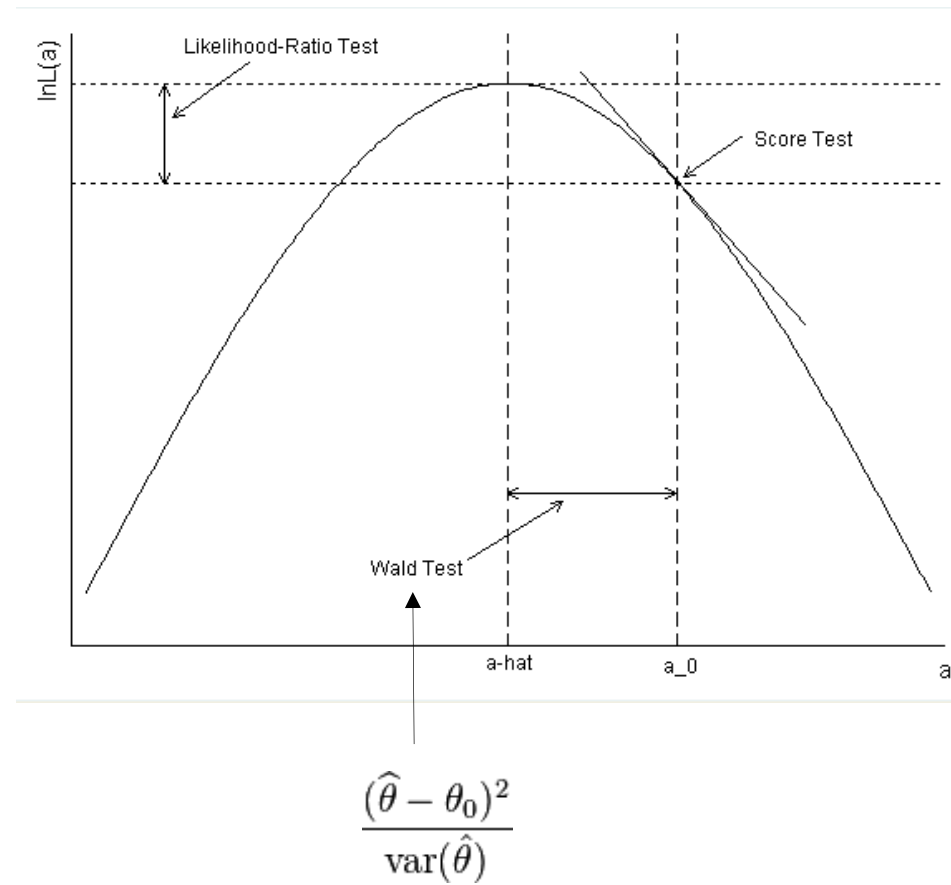
## The score statistic

- All the tests are in some sense a modification of the score test for the previous  $H_0$
- Under  $H_0$ , the score statistic for a single variant  $j$  (marginally)

$$S_j = \sum_{i=1}^n G_{ij}(Y_i - \hat{\mu}_i),$$

where  $\hat{\mu}_i$  is estimated under the null model with  $\beta$  set to zero vector

## The score test vs the Wald Test





## Burden tests


- Collapse information on multiple genetic variants into a single genetic score
- Essentially an association test between the score and trait
- Define a weight for each variant  $\mathbf{w} = (w_1, \dots, w_m)'$ , the score is developed as

$$\begin{aligned}h(\mu(Y_i)) &= \alpha_0 + \alpha' \mathbf{X}_i + \beta' \mathbf{G}_i \\ &= \alpha_0 + \alpha' \mathbf{X}_i + \tilde{\beta} \underbrace{\mathbf{w}' \mathbf{G}_i}_{\text{scalar } C_i} \quad (2)\end{aligned}$$

- Under  $H_0 : \tilde{\beta} = 0$ , the score statistic is  
 $Q_{burden} = (\sum_{j=1}^m w_j S_j)^2$

## Burden tests

- If  $\mathbf{w} = \mathbf{1}$ , we can collapse rare variants the following way

Y	G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>	G <sub>4</sub>		C
1	1	0	0	0		1
1	0	1	0	0		1
1	0	0	1	1		2
.	.	.	.	.		.
.	.	.	.	.		.
.	.	.	.	.		.
0	0	0	0	0		0
0	0	0	0	0		0
0	0	0	0	0		0

## Burden tests

- Choice of  $\mathbf{w}$  accommodates different assumptions about disease mechanism
- e.g., the cohort allelic sums test (CAST)

$$C_i = \begin{cases} 1 & \text{when } \mathbf{1}'\mathbf{G}_i > 0, \\ 0 & \text{otherwise.} \end{cases}$$

- **Limitation:** strong assumption about the same direction/magnitude of effect (post to weight adjustment); loss of power

## Key features of burden tests

- Collapse many variants into single risk score
- Several flavors exist:
  - In general they all combine rare variants into a genetic score  
Example: Combine minor allele counts into a single risk score (dominant genetic model)
  - Weighted or unweighted versions (f.i., to prioritize certain variant types, based on predictions about damaging effect)

## Some problems with burden tests

- Problem 1: When high linkage disequilibrium (LD) [allelic non-independence] exists in the “region”, combined counts may be artificially elevated
- Problem 2: Assumes that all rare variants in a set are causal and associated with a trait in the same direction
  - Counter-examples exist for different directionality (e.g. autoimmune GWAs)
  - Violations of this assumption leads to power loss

## Adaptive burden tests

- To obtain tests that are robust to null variants and allow for different effect directions
- Let the data speak!
- e.g. the data-adaptive sum test (aSum)
  - Estimate direction of each variant in [marginal models](#)
  - Use the burden test framework with  $w_j = 1$  if the coef is [likely](#) to be positive and  $w_j = -1$  otherwise
  - Require permutation ([How?](#)) to obtain the null distribution
  - Further modification based on model-selection allowing for zero weight
- [Limitations](#): although more robust, marginal models are unstable; permutation requires extensive calculation

## Variance components tests

- Is there another way to pool/group the rare variants in a region?
- Yes, resort to random-effects models
- To evaluate the distribution of genetic effects for a group of variants
- Suppose  $\beta_j \sim N(0, w_j^2 \tau)$ ,  $\text{corr}(\beta_j, \beta_k) = \rho$ 
  - e.g., the widely-used sequence kernel association test (SKAT,  $\rho = 0$ ) tests  $H_0 : \tau = 0$
  - $Q_{SKAT} = \sum_{j=1}^m w_j^2 S_j^2$ , a weighted sum of squares of single variant scores, approx follows a mixture of Chi-squared dist
  - Robust to different directions of effects, but ...
  - Can lead to inflated test size in small effective sample size

## Omnibus tests

- To achieve robust power
- Often referred to as “combined tests”
- How to combine different tests?
- Fisher’s combination method

$$Fisher = -2 \log(p_{SKAT}) - 2 \log(p_{burden})$$

- Combining test statistic

$$Q_\rho = (1 - \rho)Q_{SKAT} + \rho Q_{burden}, \quad \rho \in [0, 1]$$



## Omnibus tests

- **Limitation:** it might have lower power than SKAT or burden tests if the assumption underlying one of these tests are largely true
- For unknown genetic architecture, this is an attractive choice

## General comments on aggregation tests

- The tests are designed to boost power assuming the rare variants can be grouped together
- This point is shown by simulation work by Li and Neal, 2008
- Power loss occurs (relative to single-variant tests) when only a very few of the variants are associated with the trait and when many variants have no effects
- e.g., Liu et al. studied the association between blood lipids and BCAM and CD300LG, but found weaker signal using gene-level test than single-variant test

## Meta-analysis

- Combine data from multiple studies
  - Rare variants association detection requires large sample
- Popular frameworks combine score statistic from different studies instead of combining p-values
  - only requires summary statistics
  - allows study-specific covariates
- Methods should account for heterogeneity of genetic effects ([how?](#) see Lee et al 2013 AJHG) across studies to increase power (diff in ancestries)

## Variant selection: which variants to use?

- Can use all the variants
- Obtain a refined subset on the basis of MAF, impact of amino acid sequence
- A subset based on the predicted functional role of variants (with bioinformatics tool)

## Rare-Variant Association Analysis: Study Designs and Statistical Tests

Seunggeung Lee,<sup>1</sup> Gonçalo R. Abecasis,<sup>1</sup> Michael Boehnke,<sup>1</sup> and Xihong Lin<sup>2,\*</sup>

Despite the extensive discovery of trait- and disease-associated common variants, much of the genetic contribution to complex traits remains unexplained. Rare variants can explain additional disease risk or trait variability. An increasing number of studies are underway to identify trait- and disease-associated rare variants. In this review, we provide an overview of statistical issues in rare-variant association studies with a focus on study designs and statistical tests. We present the design and analysis pipeline of rare-variant studies and review cost-effective sequencing designs and genotyping platforms. We compare various gene- or region-based association tests, including burden tests, variance-component tests, and combined omnibus tests, in terms of their assumptions and performance. Also discussed are the related topics of meta-analysis, population-stratification adjustment, genotype imputation, follow-up studies, and heritability due to rare variants. We provide guidelines for analysis and discuss some of the challenges inherent in these studies and future research directions.

(Lee et al. 2014)

## Other tests

	<b>Description</b>	<b>Methods</b>	<b>Advantage</b>	<b>Disadvantage</b>	<b>Software Packages<sup>a</sup></b>
Burden tests	collapse rare variants into genetic scores	ARIEL test, <sup>50</sup> CAST, <sup>51</sup> CMC method, <sup>52</sup> MZ test, <sup>53</sup> WSS <sup>54</sup>	are powerful when a large proportion of variants are causal and effects are in the same direction	lose power in the presence of both trait-increasing and trait-decreasing variants or a small fraction of causal variants	EPACTS, GRANVIL, PLINK/SEQ, Rvtests, SCORE-Seq, SKAT, VAT
Adaptive burden tests	use data-adaptive weights or thresholds	aSum, <sup>55</sup> Step-up, <sup>56</sup> EREC test, <sup>57</sup> VT, <sup>58</sup> KBAC method, <sup>59</sup> RBT <sup>60</sup>	are more robust than burden tests using fixed weights or thresholds; some tests can improve result interpretation	are often computationally intensive; VT requires the same assumptions as burden tests	EPACTS, KBAC, PLINK/SEQ, Rvtests, SCORE-Seq, VAT
Variance-component tests	test variance of genetic effects	SKAT, <sup>61</sup> SSU test, <sup>62</sup> C-alpha test <sup>63</sup>	are powerful in the presence of both trait-increasing and trait-decreasing variants or a small fraction of causal variants	are less powerful than burden tests when most variants are causal and effects are in the same direction	EPACTS, PLINK/SEQ, SCORE-Seq, SKAT, VAT

(Lee et al. 2014)

## Other tests

Combined tests	combine burden and variance-component tests	SKAT-O, <sup>64</sup> Fisher method, <sup>65</sup> MiST <sup>66</sup>	are more robust with respect to the percentage of causal variants and the presence of both trait-increasing and trait-decreasing variants	can be slightly less powerful than burden or variance-component tests if their assumptions are largely held; some methods (e.g., the Fisher method) are computationally intensive	EPACTS, PLINK/SEQ, MiST, SKAT
EC test	exponentially combines score statistics	EC test <sup>67</sup>	is powerful when a very small proportion of variants are causal	is computationally intensive; is less powerful when a moderate or large proportion of variants are causal	no software is available yet

Abbreviations are as follows: ARIEL, accumulation of rare variants integrated and extended locus-specific; aSum, data-adaptive sum test; CAST, cohort allelic sums test; CMC, combined multivariate and collapsing; EC, exponential combination; EPACTS, efficient and parallelizable association container toolbox; EREC, estimated regression coefficient; GRANVIL, gene- or region-based analysis of variants of intermediate and low frequency; KBAC, kernel-based adaptive cluster; MiST, mixed-effects score test for continuous outcomes; MZ, Morris and Zeggini; RBT, replication-based test; Rvtests, rare-variant tests; SKAT, sequence kernel association test; SSU, sum of squared score; VAT, variant association tools; VT, variable threshold; and WSS, weighted-sum statistic.

<sup>a</sup>More information is given in [Table 3](#).

(Lee et al. 2014)



## A comprehensive evaluation of collapsing methods using simulated and real data: excellent annotation of functionality and large sample sizes required

Carmen Dering<sup>1</sup>, Inke R. König<sup>1</sup>, Laura B. Ramsey<sup>2</sup>, Mary V. Relling<sup>2</sup>, Wenjian Yang<sup>2</sup> and Andreas Ziegler<sup>1,3,4\*</sup>

<sup>1</sup> Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Lübeck, Germany

<sup>2</sup> Pharmaceutical Department, St. Jude Children's Research Hospital, Memphis, TN, USA

<sup>3</sup> Zentrum für Klinische Studien, Universität zu Lübeck, Lübeck, Germany

<sup>4</sup> School of Mathematics, Statistics and Computer Science, University of KwaZulu-Natal, Durban, South Africa

### Edited by:

Daniel C. Koboldt, Washington University in St. Louis, USA

### Reviewed by:

Michelle Leary, Tulane University, USA

Jian Li, Tulane University, USA

### \*Correspondence:

Andreas Ziegler, Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Campus Lübeck, Ratzeburger Allee 160, 23562 Lübeck, Germany  
e-mail: ziegler@imbs.uni-luebeck.de

The advent of next generation sequencing (NGS) technologies enabled the investigation of the rare variant-common disease hypothesis in unrelated individuals, even on the genome-wide level. Analysis of this hypothesis requires tailored statistical methods as single marker tests fail on rare variants. An entire class of statistical methods collapses rare variants from a genomic region of interest (ROI), thereby aggregating rare variants. In an extensive simulation study using data from the Genetic Analysis Workshop 17 we compared the performance of 15 collapsing methods by means of a variety of pre-defined ROIs regarding minor allele frequency thresholds and functionality. Findings of the simulation study were additionally confirmed by a real data set investigating the association between methotrexate clearance and the *SLCO1B1* gene in patients with acute lymphoblastic leukemia. Our analyses showed substantially inflated type I error levels for many of the proposed collapsing methods. Only four approaches yielded valid type I errors in all considered scenarios. None of the statistical tests was able to detect true associations over a substantial proportion of replicates in the simulated data. Detailed annotation of functionality of variants is crucial to detect true associations. These findings were confirmed in the analysis of the real data. Recent theoretical work showed that large power is achieved in gene-based analyses only if large sample sizes are available and a substantial proportion of causing rare variants is present in the gene-based analysis. Many of the investigated statistical approaches use permutation requiring high computational cost. There is a clear need for valid, powerful and fast to calculate test statistics for studies investigating rare variants.





## A comprehensive evaluation of collapsing methods using simulated and real data: excellent annotation of functionality and large sample sizes required

Carmen Dering<sup>1</sup>, Inke R. König<sup>1</sup>, Laura B. Ramsey<sup>2</sup>, Mary V. Relling<sup>2</sup>, Wenjian Yang<sup>2</sup> and Andreas Ziegler<sup>1,3,4\*</sup>

<sup>1</sup> Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Lübeck, Germany

<sup>2</sup> Pharmaceutical Department, St. Jude Children's Research Hospital, Memphis, TN, USA

<sup>3</sup> Zentrum für Klinische Studien, Universität zu Lübeck, Lübeck, Germany

<sup>4</sup> School of Mathematics, Statistics and Computer Science, University of KwaZulu-Natal, Durban, South Africa

### Edited by:

Daniel C. Koboldt, Washington University in St. Louis, USA

### Reviewed by:

Michelle Leacy, Tulane University, USA

Jian Li, Tulane University, USA

### \*Correspondence:

Andreas Ziegler, Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Lübeck, Germany  
e-mail: zieg@imbs.uni-luebeck.de

The advent of next generation sequencing (NGS) enabled the investigation of the rare variant component of the genome. However, the analysis of the genome-wide level of rare variants requires tailored statistical methods as single marker tests are not powerful enough. An entire class of statistical methods collapses rare variants into a region of interest (ROI), thereby aggregating rare variants. We evaluated the performance of 15 collapsing methods by means of a variety of simulated and real data sets regarding minor allele frequency thresholds and functionality. Findings from the simulation study were additionally confirmed by a real data set investigating the association between methotrexate clearance and the *SLCO1B1* gene in patients with acute lymphoblastic leukemia. Our analyses showed substantially inflated type I error levels for many of the proposed collapsing methods. Only four approaches yielded valid type I errors in all considered scenarios. None of the statistical tests was able to detect true associations over a substantial proportion of replicates in the simulated data. Detailed annotation of functionality of variants is crucial to detect true associations. These findings were confirmed in the analysis of the real data. Recent theoretical work showed that large power is achieved in gene-based analyses only if large sample sizes are available and a substantial proportion of causing rare variants is present in the gene-based analysis. Many of the investigated statistical approaches use permutation requiring high computational cost. There is a clear need for valid, powerful and fast to calculate test statistics for studies investigating rare variants.

Aggregation tests typically do not perform well

## Which tests to use?

- First acknowledge that relative performance depends on the unknown disease architecture
- Use available prior information
  - the region has a large fraction of causal rare variants, majority increase disease risk – burden tests
  - exist both risk-increasing and risk-decreasing variants – variance-component tests
- If no prior information, one can try multiple methods or use the omnibus test

## 5 When variants become rare – sparse data – Sequence Analyses

### Counting letters or words (see next)

- The **CpG sites** or **CG sites** are regions of DNA where a cytosine nucleotide occurs next to a guanine nucleotide in the linear sequence of bases along its length. "CpG" is shorthand for "—C—phosphate—G—", that is, cytosine and guanine separated by only one phosphate. The "CpG" notation is used to distinguish this linear sequence from the CG base-pairing of cytosine and guanine.

([https://en.wikipedia.org/wiki/CpG\\_site](https://en.wikipedia.org/wiki/CpG_site))

```

CATTCCGCTTCTCTCCGAGGTGGCGCGTGGGA
GGTGTTTTGCTCGGGTCTGTAAGAATAGGCCAGG
CAGCTTCCCGCGGGATGCGCTCATCCCCTCTCG
GGTTCCTCCACCGCGCGCGTTGGCCCGGT
CCGCTGCGAGATGTTTTCCGACGACAATGATTC
CACTCTCGCGCTCCCATGTTGATCCCAGCTCCT
CTGCGGGCGTCAGGACCCCTGGGCCCGCCCG
CTCCACTCAGTCAATCTTTTGTCCCCTATAAGCG
GATTATCGGGGTGGCTGGGGGCGGCTGATTCGA
CGAATGCCCTTGGGGGTCACC CGGAGGGAACTC
CGGGCTCCGGCTTTGGCCAGCCCGCACCCCTGGT
TGAGCCGGCCCGAGGGCCACCAGGGGGCGCTCG
ATGTTCTGCAGCCCCCGCAGCAGCCCCACTCC
CCGGCTCACCCCTACGATTGGCTGGCCCGCCCGAG
CTCTGTGCTGTGATTGGTCACAGCCCGTGTCCGTC
GCGGGCGCGGGCGGATACGAGGTGACGCGCA
GAGGCCAGCTCGGGGCGGTGTCCCGCGCGCGC
GACTGCGGGCGGAGTTTCCGCGAGGGCCGAGCG
GGGCAGTGTGACGGCAGCGGTCTGGGAGGCGC
CCGCGCGCGTCCGAGCAGCTCCCCTCCTCGCA
GCCTCACCGCGGCGTCCCGCGCCCTGGCC
TCCCGCACTCGCGCACTCCTGTCCCGCGCCACC
GCCACCTCCCACCTCGATGCGGTGC CGGGCTGC
TGCGTGATGGGGCTGCGGAGCGGCGCCCTGCGG
CTCGCGCGCGCTGCTCGCGCTGAGGTGCGT
CGGTGCCCGGCCCCCGCGCCCGCGCGCGG
CTCTTAGTTTTGGGTGCATTTGTCTGGTCTTCCAAA
CTAGATTGAAAGCTCTGAAAAAAAAAACTATCTTGT
GTTTCTATCTGTTGAGCTCATAGTAGGTATCCAGGA
AGTAGTAGGGTTGACTGCATTGATTTGGGACTACAC
TGGGAGTTTTCTTCCCATCTCCCCTTAGTTTTCT
TTTTTCTTTCTTTCTTTCTTTTTTTTTTTTTTTTT
TTGAGATGCTCTTGCTCAGTCCCCCAGGCTGGA
GTGCAGTGGTGCGATCTTGGCTCACTGTAGCCTCC
ACCTCCCAGGTTCAAGCAATCTACTGCCTTAGCCT
CCCGAGTAGCTGGGATTACAAGCACC CGCCACCAT
TCCTGGCTAATTTTTTTTTTTGTATTTTAGTTGAGA
CAGGGTTTACCATGTTGGTGTGCTGGTCTCAGA
CTCCTGGGGCCTAGCGATCCCCCTGCCTCAGCCT
CCCAGAGTGTAGGATTACAGGCATGAGCCACTGT
ACC CGCCTCTCTCCAGTTTCCAGTTGGAATCCAA
GGGAAGTAAGTTTAAGATAAAGTTACGATTTTGAAT
CTTTGGATTGAGAAGAATTTGTCACCTTTAACACCT
AGAGTTGAACTTCATACCTGGAGAGCCTTAACATT
AAGCCCTAGCCAGCCTCCAGCAAGTGGACATTGGT
CAGGTTTGGCAGGATTCTCCCCTGAAGTGGACT
GAGAGCCACACCCTGGCCTGTACCATACCCATCC
CCTATCCTTAGTGAAGCAAACTCCTTTGTTCCCTT
CTCCTTCTCCTAGTGACAGGAAATATTGTGATCCTA
AAGAATGAAAATAGCTTGTACCTCGTGGCCTCAG
GCCCTTGACTTCAGGCGTTCTGTTTAAATCAAGT
GACATCTTCCCGAGGCTCCCTGAATGTGGCAGATG

```

## Comparing multiple sequences (see practical session)

- After collection of a set of related sequences, how can we compare them as a set?
- How should we line up the sequences so that the most similar portions are together?
- What do we do with sequences of different length?

```

                2430          2440          2450          2460          2470
HSA128 CACTTCCCCTAT---GCAGGTGTCCAACGGATGTGTGAGTAAAATTCTGGGCAGGTATTA
      ::  :::::  ::  ::::::::::::::::::::::::::::::::::::::::::::::::::::
pax6  CATTTCCCGAATTCTGCAGGTGTCCAACGGATGTGTGAGTAAAATTCTGGGCAGGTATTA
      540          550          560          570          580          590

                2480          2490          2500          2510          2520          2530
HSA128 CGAGACTGGCTCCATCAGACCCAGGGCAATCGGTGGTAGTAAACCGAGAGTAGCGACTCC
      ::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
pax6  CGAGACTGGCTCCATCAGACCCAGGGCAATCGGTGGTAGTAAACCGAGAGTAGCGACTCC
      600          610          620          630          640          650

                2540          2550          2560          2570          2580          2590
HSA128 AGAAGTTGTAAGCAAAATAGCCCAGTATAAGCGGGAGTGCCCGTCCATCTTTGCTTGGA

```

# Investigating frequencies of occurrences of words

## Introduction

- Words are short strings of letters drawn from an alphabet
- In the case of DNA, the set of letters is A, C, T, G
- A word of length  $k$  is called a  $k$ -word or  $k$ -tuple
- Differences in word frequencies help to differentiate between different DNA sequence sources or regions
- Examples: 1-tuple: individual nucleotide; 2-tuple: dinucleotide; 3-tuple: codon
- The **distributions** of the nucleotides over the DNA sequences have been studied for many years → hidden correlations in the sequences (e.g., CpGs)

## Probability distributions

### Probability is the science of uncertainty

1. Rules → data: given the rules, describe the likelihoods of various events occurring
2. Probability is about prediction – looking forwards
3. Probability is mathematics

## Statistics is the science of data

1. Rules  $\leftarrow$  data: given only the data, try to guess what the rules were. That is, some probability model controlled what data came out, and the best we can do is guess – or approximate – what that model was. We might guess wrong, we might refine our guess as we obtain / collect more data
2. Statistics is about looking backward. Once we make our best *statistical guess* about what the probability model is (what the rules are), based on looking backward, we can then use that probability model to predict the future
3. Statistics is an art. It uses mathematical methods but it is much more than mathematics alone
4. The purpose of statistics is to make inference about unknown quantities from samples of data.



## Statistics is the science of data

- Probability distributions are a fundamental concept in statistics.
- Before computing an interval or test based on a distributional assumption, we need to verify that the assumption is justified for the given data set.
- For this chapter, the distribution does not always need to be the best-fitting distribution for the data, but an adequate enough model so that the statistical technique yields valid conclusions.
- Simulation studies: one way to obtain empirical evidence for a probability model

## Expected values and variances

- Mean and variance are two important properties of real-valued random variables and corresponding probability distributions.
- The “mean” of a discrete random variable  $X$  taking values  $x_1, x_2, \dots$  (denoted  $EX$  (or  $E(X)$  or  $E[X]$ ), where  $E$  stands for expectation, which is another term for mean) is defined as:

$$E(X) = \sum_i x_i P(X = x_i)$$

- $E(X_i) = 1 \times p_A + 0 \times (1 - p_A)$  if  $x_i = A$  or {another letter}
  - If  $Y = c X$ , then  $E(Y) = c E(X)$
  - $E(X_1 + \dots + X_n) = E(X_1) + \dots + E(X_n)$
- Because  $X_i$  are assumed to be independent and identically distributed (iid):

$$E(X_1 + \dots + X_n) = n E(X_1) = n p_A$$

## Expected values and variances

- The idea is to use squared deviations of  $X$  from its center (expressed by the mean). Expanding the square and using the linearity properties of the mean, the  $\text{Var}(X)$  can also be written as:

$$\text{Var}(X) = E(X^2) - [E(X)]^2$$

- If  $Y=c X$  then  $\text{Var}(Y) = c^2 \text{Var}(X)$
  - The variance of a sum of independent random variables is the sum of the individual variances
- 
- For the random variables  $X_i$  taking on values  $A$  or sth else:  
 $\text{Var}(X_i) = [1^2 \times p_A + 0^2 \times (1 - p_A)] - p_A^2 = p_A(1 - p_A)$   
 $\text{Var}(N) = n \text{Var}(X_1) = np_A(1 - p_A)$

## Expected values and variances

- The expected value of a random variable  $X$  gives a measure of its location. Variance is another property of a probability distribution dealing with the spread or variability of a random variable around its mean.

$$\text{Var}(X) = E ( [X - E(X)]^2 )$$

- The positive square root of the variance of  $X$  is called its standard deviation  $\text{sd}(X)$

## Independence

- Discrete random variables  $X_1, \dots, X_n$  are said to **be independent** if for any subset of random variables and actual values, the joint distribution equals the product of the component distributions

## Occurrences of 1-letter words

### Assumptions

- Notation for the output of a random string of  $n$  bases may be:  $L_1, L_2, \dots, L_n$   
( $L_i$  = base inserted at position  $i$  of the sequence)
  - The values  $l_j$  for  $L_j$  will come from a set  $\chi$  (with  $J$  possibilities)
  - For a DNA sequence,  $J=4$  and  $\chi = \{A, C, T, G\}$
- Simple rules specifying a probability model:
  - First base in sequence is either A, C, T or G with prob  $p_A, p_C, p_T, p_G$
  - Suppose the first  $r$  bases have been generated, while generating the base at position  $r+1$ , no attention is paid to what has been generated before.

- Then we can actually generate A, C, T or G with the probabilities above
- According to our simple model, the  $L_i$  are independent and hence

$$P(L_1=l_1, L_2=l_2, \dots, L_n=l_n) = P(L_1=l_1) P(L_2=l_2) \dots P(L_n=l_n)$$

- If  $p_j$  is the prob that the value (realization of the random variable  $L$ )  $l_j$  occurs, then

- $p_1, \dots, p_J \geq 0$  and  $p_1 + \dots + p_J = 1$

- The **probability distribution** (probability mass function) of  $L$  is given by the collection  $p_1, \dots, p_J$

- $P(L=l_j) = p_j, j=1, \dots, J$

- The probability that an event  $S$  occurs (subset of  $\chi$ ) is  $P(L \in S) = \sum_{j:l_j \in S} (p_j)$

## Probability distributions of interest

- What is the probability distribution of the number of times a given pattern occurs in a random DNA sequence  $L_1, \dots, L_n$ ?

- New sequence  $X_1, \dots, X_n$ :

$$X_i=1 \text{ if } L_i=A \text{ and } X_i=0 \text{ else}$$

- The number of times  $N$  that  $A$  appears is the sum

$$N=X_1+\dots+X_n$$

- The prob distr of each of the  $X_i$ :

$$P(X_i=1) = P(L_i=A)=p_A$$

$$P(X_i=0) = P(L_i=C \text{ or } G \text{ or } T) = 1 - p_A$$

- What is a “typical” value of  $N$ ?

- Depends on how the individual  $X_i$  (for different  $i$ ) are interrelated



## The binomial distribution

- The binomial distribution is used when there are exactly two mutually exclusive outcomes of a trial. These outcomes are appropriately labeled "success" and "failure". The binomial distribution is used to obtain the probability of observing  $x$  successes in a fixed number of trials, with the probability of success on a single trial denoted by  $p$ . The binomial distribution assumes that  $p$  is fixed for all trials.
- The formula for the binomial probability mass function is :

$$P(N = j) = \binom{n}{j} p^j (1 - p)^{n-j}, j = 0, 1, \dots, n$$

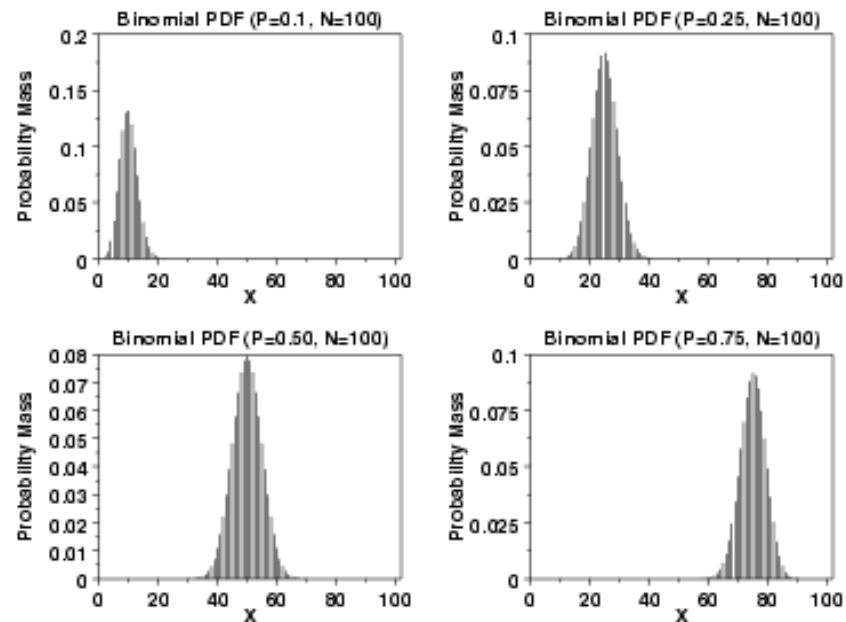
with the binomial coefficient  $\binom{n}{j}$  determined by

$$\binom{n}{j} = \frac{n!}{j! (n - j)!}$$

and  $j! = j(j-1)(j-2)\dots 3.2.1$ ,  $0! = 1$

## The binomial distribution

- The mean is  $np$  and the variance is  $np(1-p)$
- The following is the plot of the binomial probability density function for four values of  $p$  and  $n = 100$ .



## Simulating from probability distributions

- The idea is that we can study the properties of the distribution of  $N$  when we can get our computer to output numbers  $N_1, \dots, N_n$  having the same distribution as  $N$

- We can use the sample mean to estimate the expected value  $E(N)$ :

$$\bar{N} = (N_1 + \dots + N_n)/n$$

- Similarly, we can use the sample variance to estimate the true variance of  $N$ :

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (N_i - \bar{N})^2$$

Why do we use  $(n-1)$  and not  $n$  in the denominator?

## Simulating from probability distributions

- What is needed to produce such a string of observations?
  - Access to **pseudo-random numbers**: random variables that are uniformly distributed on (0,1): any number between 0 and 1 is a possible outcome and each is equally likely
- In practice, simulating an observation with the distribution of  $X_1$ :
  - Take a uniform random number  $u$
  - Set  $X_1=1$  if  $U \leq p \equiv p_A$  and 0 otherwise.
  - Why does this work? ...  $P(X_1 = 1) = P(U \leq p_A) = p_A$
  - Repeating this procedure  $n$  times results in a sequence  $X_1, \dots, X_n$  from which  $N$  can be computed by adding the  $X$ 's

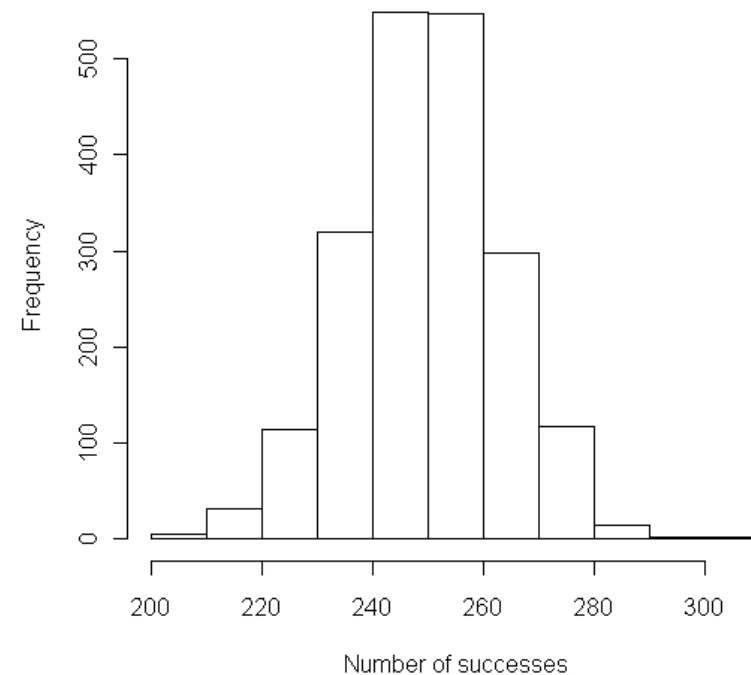
## Simulating from probability distributions

- Simulate a general DNA sequence of bases  $L_1, \dots, L_n$ :
  - Divide the interval  $(0,1)$  in 4 intervals with endpoints
$$0, p_A, p_A + p_C, p_A + p_C + p_G, 1$$
  - If the simulated  $u$  lies in the leftmost interval,  $L_1=A$
  - If  $u$  lies in the second interval,  $L_1=C$ ; if in the third,  $L_1=G$  and otherwise  $L_1=T$
  - Repeating this procedure  $n$  times with different values for  $U$  results in a sequence  $L_1, \dots, L_n$
- Use the “sample” function in R:

```
pi <- c(0.25,0.75)
x<-c(1,0)
set.seed(2009)
sample(x,10,replace=TRUE,pi)
```

## Simulating from probability distributions

- By looking through a given simulated sequence, we can count the number of times a particular pattern arises (for instance, the base A)
- By repeatedly generating sequences and analyzing each of them, we can get a feel for whether or not our particular pattern of interest is unusual



```
x<- rbinom(2000,1000,0.25)
mean(x)
sd(x)^2
hist(x,xlab="Number of successes",main="")
```

# R documentation

Binomial {stats}

R Documentation

## The Binomial Distribution

### Description

Density, distribution function, quantile function and random generation for the binomial distribution with parameters `size` and `prob`.

This is conventionally interpreted as the number of ‘successes’ in `size` trials.

### Usage

```
dbinom(x, size, prob, log = FALSE)
pbinom(q, size, prob, lower.tail = TRUE, log.p = FALSE)
qbinom(p, size, prob, lower.tail = TRUE, log.p = FALSE)
rbinom(n, size, prob)
```

### Arguments

<code>x, q</code>	vector of quantiles.
<code>p</code>	vector of probabilities.
<code>n</code>	number of observations. If <code>length(n) &gt; 1</code> , the length is taken to be the number required.
<code>size</code>	number of trials (zero or more).

<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/Binomial.html>

```
> rbinom(1,1000,0.25)
```

```
[1] 250 → you got lucky!!!!
```

## Simulating from probability distributions

- Using R code:

```
x<- rbinom(2000,1000,0.25)
mean(x)
sd(x)^2
hist(x,xlab="Number of successes",main="")
```

What is the number of observations?



## Simulating from probability distributions

- Using R code:

```
x<- rbinom(2000,1000,0.25)
mean(x)
sd(x)^2
hist(x,xlab="Number of successes",main="")
```

What is the number of observations?

Number of sequences = 2000  
Number of trials = 1000

## Back to our original question

- Suppose we have a sequence of 1000bp and assume that every base occurs with equal probability. How likely are we to observe at least 300 A's in such a sequence?
  - Exact computation using a closed form of the relevant distribution
  - Approximate via simulation
  - Approximate using the Central Limit Theory

## Exact computation via closed form of relevant distribution

- The formula for the binomial probability mass function is :

$$P(N = j) = \binom{n}{j} p^j (1 - p)^{n-j}, j = 0, 1, \dots, n$$

and therefore

$$\begin{aligned} P(N \geq 300) &= \sum_{j=300}^{1000} \binom{1000}{j} (1/4)^j (1 - 1/4)^{1000-j} \\ &= 0.00019359032194965841 \end{aligned}$$

- Note that the probability  $P(N \geq 300)$  is estimated to be 0.0001479292 via

`1-pbinom(300,size=1000,prob=0.25)`

`pbinom(300,size=1000,prob=0.25,lower.tail=FALSE)`

	P: exactly 300 out of 1000	
Method 1. exact binomial calculation	0.00004566114740576488	
Method 2. approximation via normal	0.000038	
Method 3. approximation via Poisson	-----	
	P: 300 or fewer out of 1000	
Method 1. exact binomial calculation	0.9998520708293378	
Method 2. approximation via normal	0.999885	
Method 3. approximation via Poisson	-----	
	P: 300 or more out of 1000	
Method 1. exact binomial calculation	0.00019359032194965841	
Method 2. approximation via normal	0.000153	
Method 3. approximation via Poisson	-----	
For hypothesis testing		
	One-Tail	Two-Tail
Method 1. exact binomial calculation	0.00019359032194965841	0.0003025705168772097
Method 2. approximation via normal	0.000153	0.000306
Method 3. approximation via Poisson	-----	-----

(<http://faculty.vassar.edu/lowry/binomialX.html>)

## Approximate via simulation

- Using R code and simulations from the theoretical distribution,  $P(N \geq 300)$  can be estimated as 0.000196 via

```
x<- rbinom(1000000,1000,0.25)
sum(x>=300)/1000000
```

## Approximate via Central Limit Theory

- The central limit theorem offers a 3<sup>rd</sup> way to compute probabilities of a distribution
- It applies to sums or averages of iid random variables
- Assuming that  $X_1, \dots, X_n$  are iid random variables with mean  $\mu$  and variance  $\sigma^2$ , then we know that for the sample average

$$\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n),$$

$$E(\bar{X}_n) = \mu \text{ and } \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

- Hence,

$$E\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right) = 0, \text{Var}\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right) = 1$$

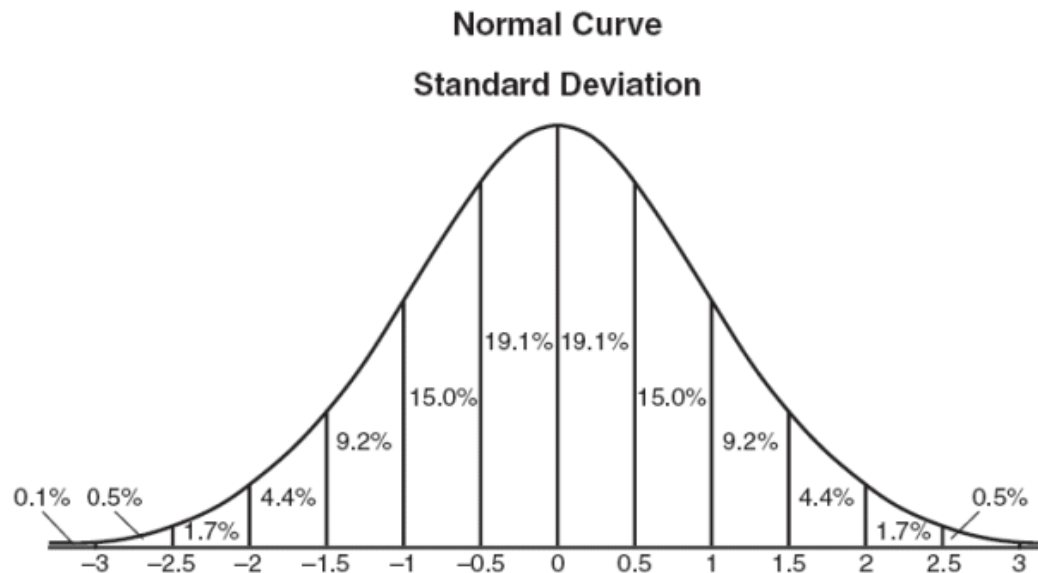
## Approximate via Central Limit Theory

- The central limit theorem states that if the sample size  $n$  is large enough,

$$P\left(a \leq \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \leq b\right) \approx \phi(b) - \phi(a),$$

with  $\phi(\cdot)$  the standard normal distribution defined as

$$\phi(z) = P(Z \leq z) = \int_{-\infty}^z \phi(x) dx$$



## Approximate via Central Limit Theory

- Estimating the quantity  $P(N \geq 300)$  when  $N$  has a binomial distribution with parameters  $n=1000$  and  $p=0.25$ ,

$$E(N) = n\mu = 1000 \times 0.25 = 250,$$

$$sd(N) = \sqrt{n} \sigma = \sqrt{1000 \times \frac{1}{4} \times \frac{3}{4}} \approx 13.693$$

$$P(N \geq 300) = P\left(\frac{N - 250}{13.693} > \frac{300 - 250}{13.693}\right)$$

$$\approx P(Z > 3.651501) = 0.0001303560$$

- R code:

```
pnorm(3.651501,lower.tail=FALSE)
```

How do the estimates of  $P(N \geq 300)$  compare?

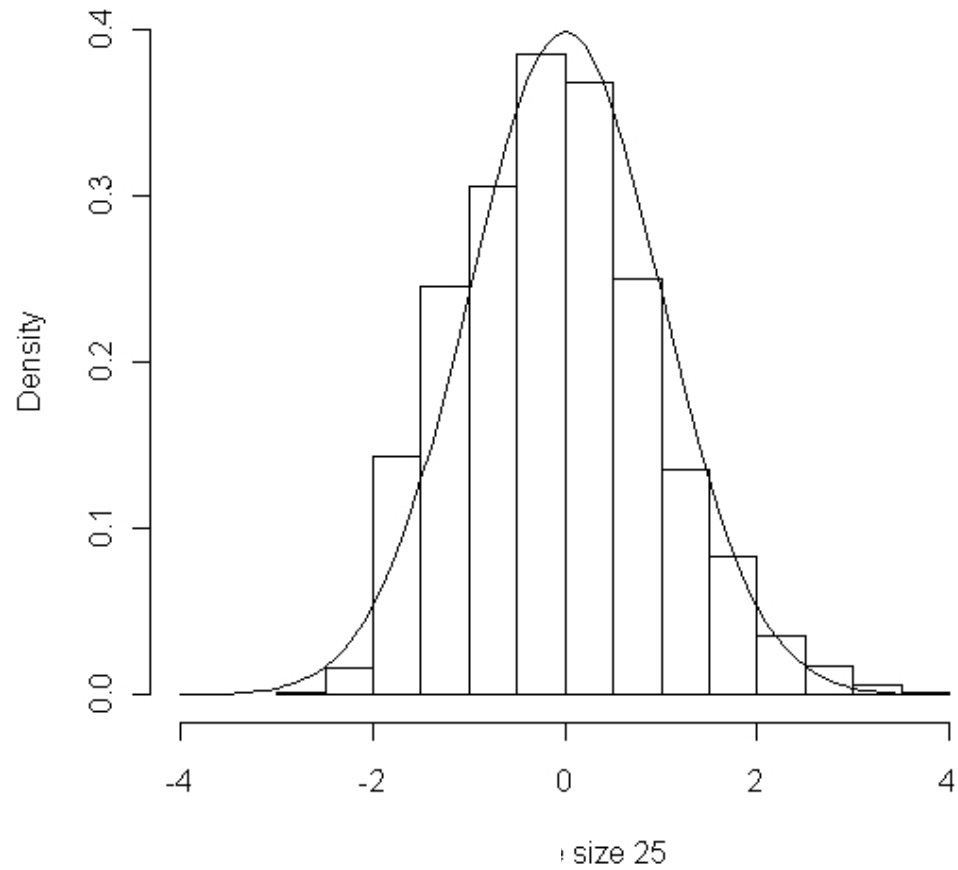


## Approximate via Central Limit Theory

- The central limit theorem in action using R code:

```
bin25<-rbinom(1000,25,0.25)
av.bin25 <- 25*0.25
stdev.bin25 <- sqrt(25*0.25*0.75)
bin25<-(bin25-av.bin25)/stdev.bin25
hist(bin25,xlim=c(-4,4),ylim=c(0.0,0.4),prob=TRUE,xlab="Sample size
25",main="")
x<-seq(-4,4,0.1)
lines(x,dnorm(x))
```

## Approximate via Central Limit Theory



## Occurrences of 2-letter words

- Concentrating on abundances, and assuming the iid model for  $L_1, \dots, L_n$ :

$$P(L_i = l_i = C, L_{i+1} = l_{i+1} = G) = p_{l_i} p_{l_{i+1}}$$

- Has a given sequence an unusual dinucleotide frequency compared to the iid model?
  - Compare observed  $O$  with expected  $E$  dinucleotide numbers

$$\chi^2 = \frac{(O-E)^2}{E},$$

with  $E = (n - 1)p_{l_i}p_{l_{i+1}}$ .

Where have we seen this statistic before? How many df?

Why  $(n-1)$  as factor in  $E$  above? How many df?

## Comparing to the reference

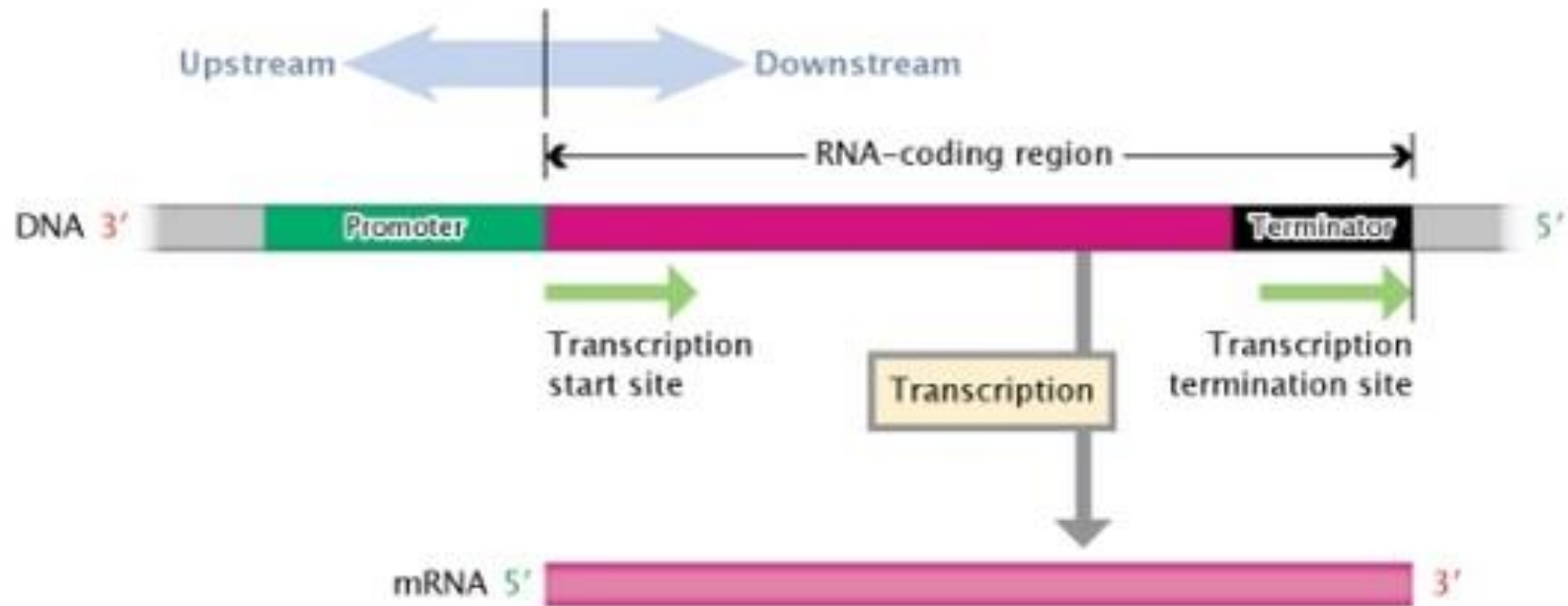
- How to determine which values of  $\chi^2$  are unlikely or extreme?
  - If the observed nr is close to the expected number, then the statistic will be small. Otherwise, the model will be doing a poor job of predicting the dinucleotide frequencies and the statistic will tend to be large...
  - Recipe:
    - Compute the number  $c$  given by
$$c = \begin{cases} 1 + 2p_{l_i} - 3p_{l_i}^2, & \text{if } l_i = l_{i+1} \\ 1 - 3p_{l_i}p_{l_{i+1}}, & \text{if } l_i \neq l_{i+1} \end{cases}$$
    - Calculate the ratio  $\frac{\chi^2}{c}$ , where  $\chi^2$  is given as before
    - If this ratio is larger than 3.84 then conclude that the iid model is not a good fit. Note that  $qchisq(0.95,1) = 3.84$

## Occurrences of 3-letter words

### Amino acids

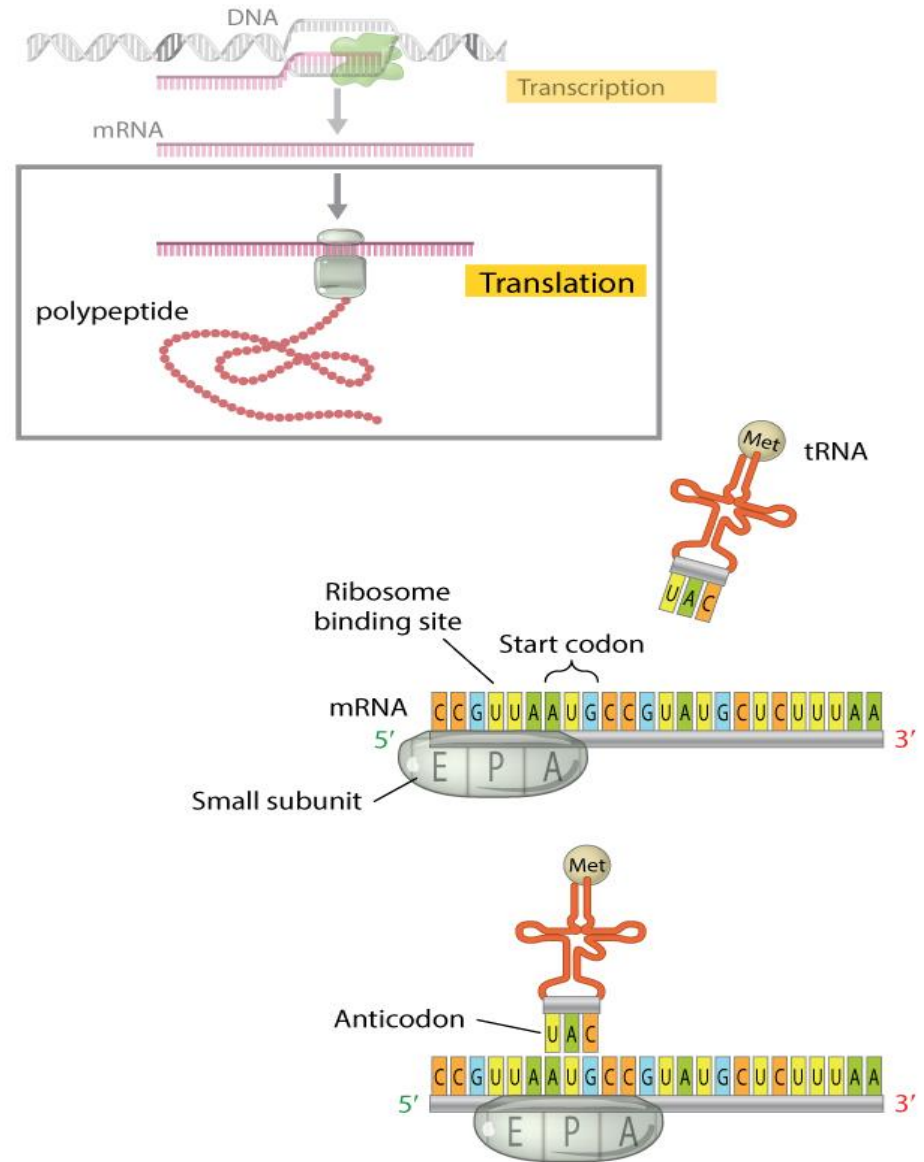
- There are 61 codons that specify amino acids and three stop codons → 64 meaningful 3-words.
- Since there are 20 common amino acids, this means that most amino acids are specified by more than one codon.

# Translation



(<https://www.nature.com/scitable>)

# Translation and Transcription

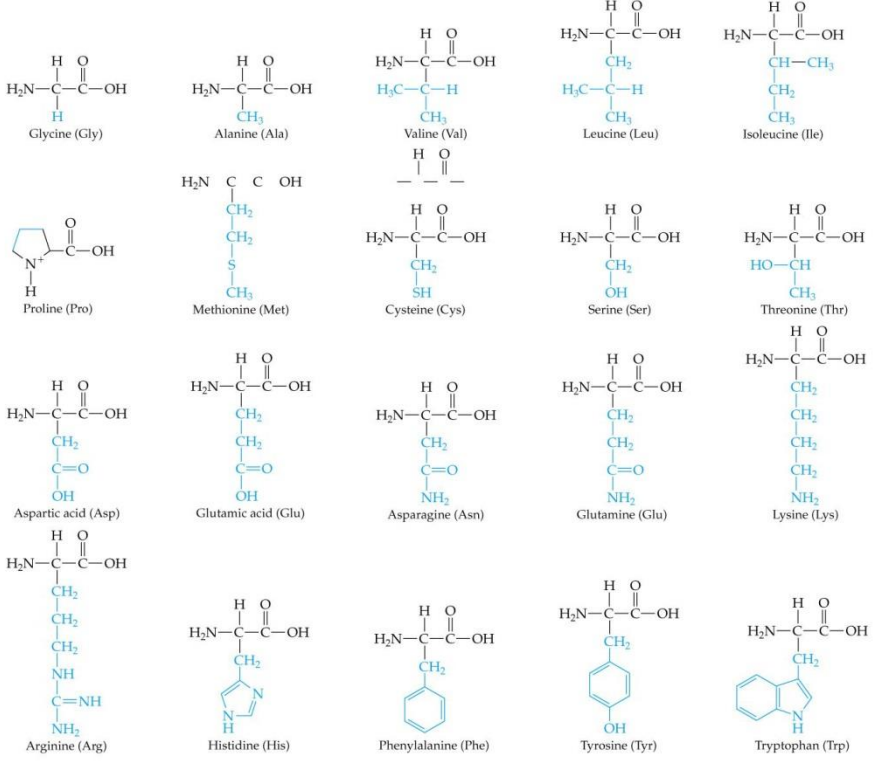


(<https://www.nature.com/scitable>)

# Amino acids

		2nd base in codon						
		U	C	A	G			
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr <b>STOP</b> <b>STOP</b>	Cys Cys <b>STOP</b> Trp	U C A G		
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G		
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G		
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G		

3rd base in codon



- This has led to the use of a number of statistics to summarize the "bias" in codon usage: An amino acid may be coded in different ways, but perhaps some codes have a preference? (higher frequency?)



## Predicted relative frequencies

- For a sequence of independent bases  $L_1, L_2, \dots, L_n$  the expected 3-tuple relative frequencies can be found by using the logic employed for dinucleotides we derived before
- The probability of a 3-word can be calculated as follows:

$$\mathbb{P}(L_i = r_1, L_{i+1} = r_2, L_{i+2} = r_3) = \mathbb{P}(L_i = r_1)\mathbb{P}(L_{i+1} = r_2)\mathbb{P}(L_{i+2} = r_3).$$

assuming the iid model

- This provides the expected frequencies of particular codons, using the individual base frequencies. It follows that among those codons making up the amino acid Phe, the expected proportion of TTT is

$$\frac{P(\text{TTT})}{P(\text{TTT}) + P(\text{TTC})}$$

## The codon adaptation index

- One can then compare predicted and observed triplet frequencies in coding sequences for a subset of genes and codons from *E. coli*.
- Médigue et al. (1991) clustered different genes based on codon usage patterns, and they observed three classes.
- For instance for Phe, the observed frequency differs considerably from the predicted frequency, when focusing on highly expressed genes (so-called “class II genes” in the work of Médigue et al. (1999) - see also next slide
- Checking the gene annotations for class II genes: highly expressed genes (ribosomal proteins or translation factors)

- Table 2.3 from Deonier et al 2005: figures in parentheses below each gene class show the number of genes in that class.

			Observed	
Codon Predicted			Gene Class I (502)	Gene Class II (191)
Phe	TTT	0.493	0.551	0.291
	TTC	0.507	0.449	0.709
Ala	GCT	0.246	0.145	0.275
	GCC	0.254	0.276	0.164
	GCA	0.246	0.196	0.240
	GCG	0.254	0.382	0.323
Asn	AAT	0.493	0.409	0.172
	AAC	0.507	0.591	0.828

**Class II : Highly expressed genes**

Class I : Moderately expressed genes

Main reference of foregoing material in this chapter: Deonier et al. *Computational Genome Analysis*, 2005, Springer (Ch 6,7)

## Supporting doc to this class (complementing course slides)



---

### REVIEW

## Rare-Variant Association Analysis: Study Designs and Statistical Tests

Seunggeung Lee,<sup>1</sup> Gonçalo R. Abecasis,<sup>1</sup> Michael Boehnke,<sup>1</sup> and Xihong Lin<sup>2,\*</sup>

Despite the extensive discovery of trait- and disease-associated common variants, much of the genetic contribution to complex traits remains unexplained. Rare variants can explain additional disease risk or trait variability. An increasing number of studies are underway to identify trait- and disease-associated rare variants. In this review, we provide an overview of statistical issues in rare-variant association studies with a focus on study designs and statistical tests. We present the design and analysis pipeline of rare-variant studies and review cost-effective sequencing designs and genotyping platforms. We compare various gene- or region-based association tests, including burden tests, variance-component tests, and combined omnibus tests, in terms of their assumptions and performance. Also discussed are the related topics of meta-analysis, population-stratification adjustment, genotype imputation, follow-up studies, and heritability due to rare variants. We provide guidelines for analysis and discuss some of the challenges inherent in these studies and future research directions.

AJHG 2014; 95, 5-23

**Questions?**