MSc GBE Course:
*Genes: from sequence to function*

**Genome-wide Association Studies**

*Sven Bergmann*

Department of Medical Genetics
University of Lausanne
Rue de Bugnon 27 - DGM 328
CH-1005 Lausanne
Switzerland

work: ++41-21-692-5452
cell: ++41-78-663-4980
http://serverdgm.unil.ch/bergmann

---

# Overview

- Population stratification
- Associations: Basics
- Whole genome associations
- Genotype imputation
- Uncertain genotypes
- New Methods

---

# Overview

- Population stratification
- Associations: Basics
- Whole genome associations
- Genotype imputation
- Uncertain genotypes
- New Methods

---

## CoLaus = Cohort Lausanne



6'189 individuals

Genotypes    Phenotypes

500.000 SNPs    159 measurement
144 questions

**Collaboration with:**
**Vincent Mooser (GSK),** *Peter Vollenweider & Gerard Waeber (CHUV)*

---

## Genetic variation in SNPs
## (**S**ingle **N**ucleotide **P**olymorphisms)



ATTGCAA**T**CCGTGG...ATC**G**AGCCA...TACGATTGCA**C**GCCG...

ATTGCAA**G**CCGTGG...ATC**T**AGCCA...TACGATTGCA**A**GCCG...

ATTGCAA**G**CCGTGG...ATC**T**AGCCA...TACGATTGCA**A**GCCG...

ATTGCAA**T**CCGTGG...ATC**G**AGCCA...TACGATTGCA**C**GCCG...

ATTGCAA**G**CCGTGG...ATC**T**AGCCA...TACGATTGCA**A**GCCG...

---

## Analysis of Genotypes only



Principle Component Analysis reveals SNP-vectors
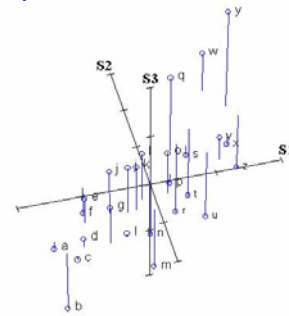explaining largest variation in the data

## Example: 2PCs for 3d-data



Raw data points: {a, …, z}

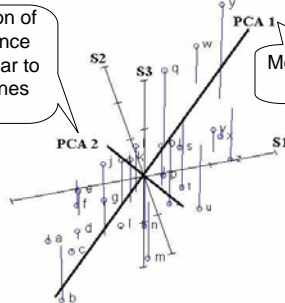## Example: 2PCs for 3d-data



Normalized data points: zero mean (& unit std)!

## Example: 2PCs for 3d-data

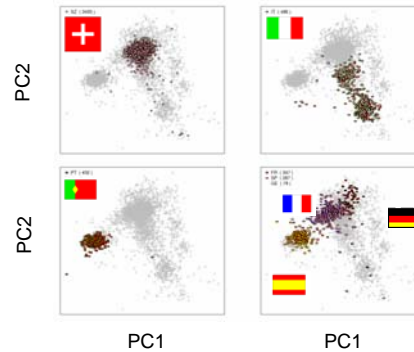The direction of most variance perpendicular to PCA1 defines PCA2
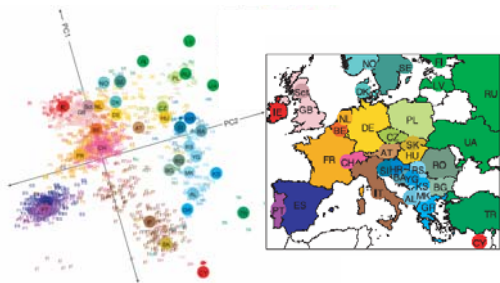
Most variance is along PCA1



Identification of axes with the most variance

## Ethnic groups cluster according to geographic distances



PC2

PC2

PC1          PC1

## PCA of POPRES cohort



**Genes mirror geography within Europe**

John Novembre[1,2], Toby Johnson[4,5,6], Katarzyna Bryc[7], Zoltán Kutalik[4,6], Adam R. Boyko[7], Adam Auton[7], Amit Indap[7], Karen S. King[8], Sven Bergmann[4,6], Matthew R. Nelson[8], Matthew Stephens[2,3] & Carlos D. Bustamante[7]

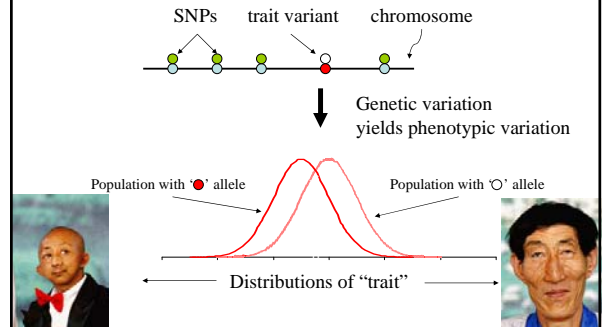nature          Vol 456|6 November 2008|doi:10.1038/nature07331

## Overview

- Population stratification
- **Associations: Basics**
- Whole genome associations
- Genotype imputation
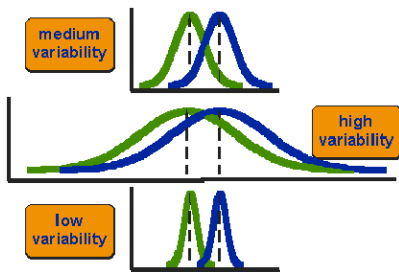- Uncertain genotypes
- New Methods

# Phenotypic variation:



# What is association?

SNPs  trait variant  chromosome

Genetic variation yields phenotypic variation

Population with '●' allele                    Population with '○' allele

Distributions of "trait"

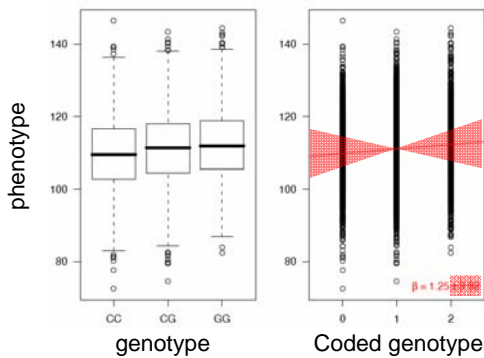# Quantifying Significance



medium variability

high variability

low variability

# T-test

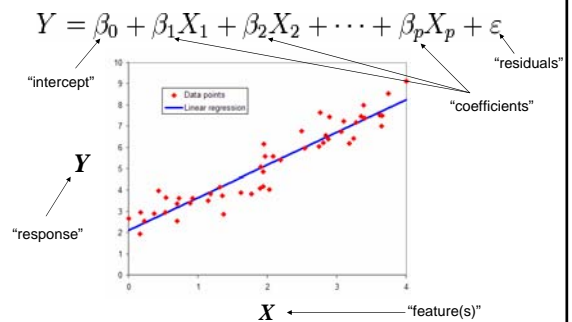$$\frac{\text{signal}}{\text{noise}} = \frac{\text{difference between group means}}{\text{variability of groups}}$$

$$\frac{\bar{x}_T - \bar{x}_C}{\sqrt{\frac{var_T}{n_T} + \frac{var_C}{n_C}}}$$
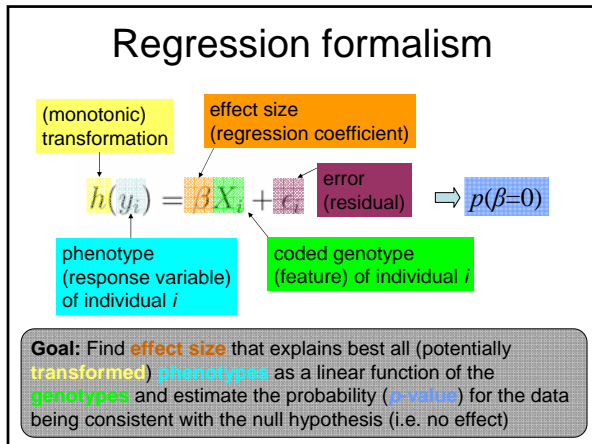
t-value

*t*-value (significance) can be translated into *p*-value (probability)
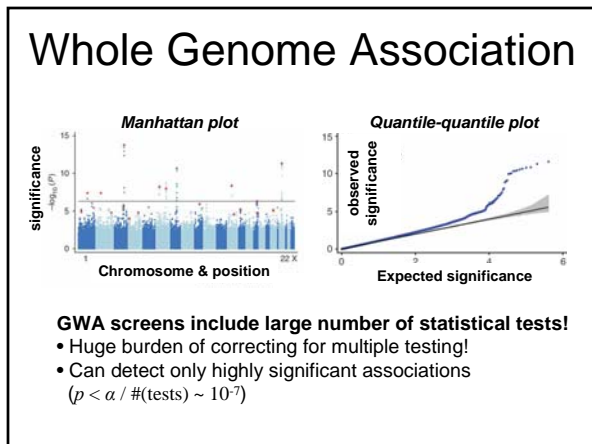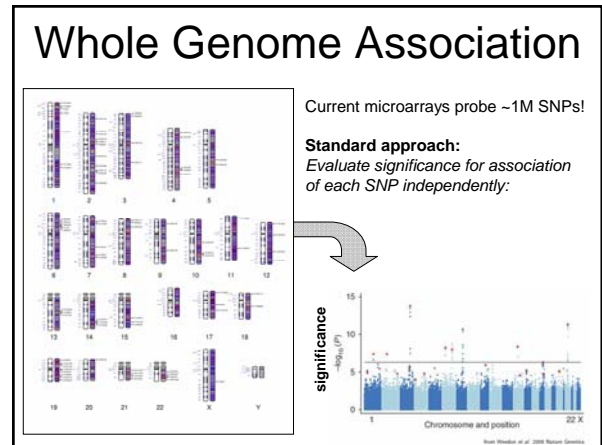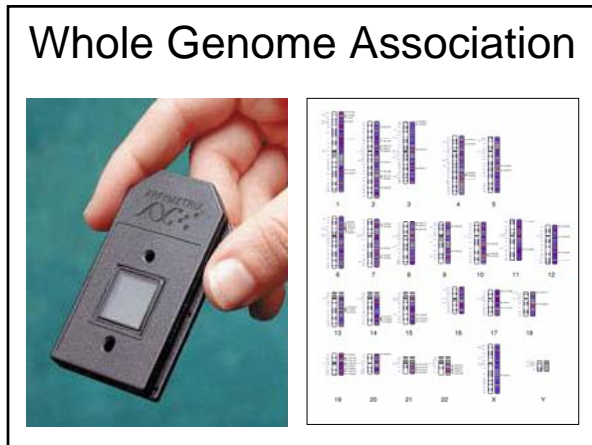
# Association using regression



phenotype

CC    CG    GG
genotype

0    1    2
Coded genotype

β = 1.25

# Regression analysis

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \varepsilon$$

"intercept"

"residuals"

"coefficients"

*Y*

"response"

Data points
Linear regression

*X*    "feature(s)"

## Regression formalism

(monotonic) transformation

effect size (regression coefficient)

error (residual)

$$h(y_i) = \beta X_i + \epsilon \implies p(\beta = 0)$$

phenotype (response variable) of individual $i$

coded genotype (feature) of individual $i$

**Goal:** Find **effect size** that explains best all (potentially **transformed**) **phenotypes** as a linear function of the **genotypes** and estimate the probability (**p-value**) for the data being consistent with the null hypothesis (i.e. no effect)

---

## Overview

- Population stratification
- Associations: Basics
- **Whole genome associations**
- Genotype imputation
- Uncertain genotypes
- New Methods

---

## Whole Genome Association



---

## Whole Genome Association



Current microarrays probe ~1M SNPs!

**Standard approach:**
*Evaluate significance for association of each SNP independently:*

significance

---

## Whole Genome Association

*Manhattan plot*

significance

Chromosome & position

*Quantile-quantile plot*

observed significance

Expected significance

**GWA screens include large number of statistical tests!**
- Huge burden of correcting for multiple testing!
- Can detect only highly significant associations
  $(p < \alpha\ /\ \#(\text{tests}) \sim 10^{-7})$

---

## GWAS: >20 publications in 2006/2007



*Massive!*

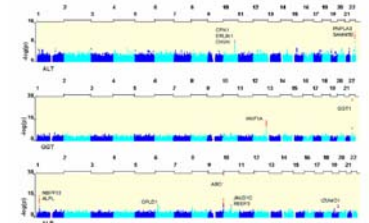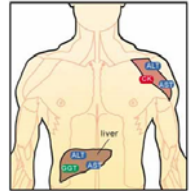## Genome-wide association analysis identifies 20 loci that influence adult height

Michael N Weedon[1,2,23], Hana Lango[1,2,23], Cecilia M Lindgren[3,4], Chris Wallace[5], David M Evans[6], Massimo Mangino[7], Rachel M Freathy[1,2], John R B Perry[1,2], Suzanne Stevens[7], Alistair S Hall[8], Nilesh J Samani[7], Beverly Shields[2], Inga Prokopenko[3,4], Martin Farrall[9], Anna Dominiczak[10], Diabetes Genetics Initiative[21], The Wellcome Trust Case Control Consortium[21], Toby Johnson[11–13], Sven Bergmann[11,12], Jacques S Beckmann[11,14], Peter Vollenweider[15], Dawn M Waterworth[16], Vincent Mooser[16], Colin N A Palmer[17], Andrew D Morris[18], Willem H Ouwehand[19,20], Cambridge GEM Consortium[22], Mark Caulfield[5], Patricia B Munroe[5], Andrew T Hattersley[1,2] & Timothy M Frayling[1,2]

## Population-Based Genome-wide Association Studies Reveal Six Loci Influencing Plasma Levels of Liver Enzymes

Xin Yuan,[1] Dawn Waterworth,[1] John R.B. Perry,[2] Noha Lim,[1] Kijoung Song,[1] John C. Chambers,[4] Weihua Zhang,[4] Peter Vollenweider,[5] Heide Stirnadel,[2] Toby Johnson,[6,7,8] Sven Bergmann,[6,8] Noam D. Beckmann,[6] Yun Li,[12] Luigi Ferrucci,[9] David Melzer,[10] Dena Hernandez,[10] Andrew Singleton,[10] James Scott,[11] Paul Elliott,[4] Gerard Waeber,[5] Lon Cardon,[3] Timothy M. Frayling,[2] Jaspal S. Kooner,[11] and Vincent Mooser[1,*]

## Common variants near *MC4R* are associated with fat mass, weight and risk of obesity

Ruth J F Loos[*,1,2,73], Cecilia M Lindgren[3,4,73], Shengxu Li[1,2], Eleanor Wheeler[5], Jing Hua Zhao[1,2], Inga Prokopenko[3,4], Michael Inouye[5], Rachel M Freathy[6,7], Antony P Attwood[5,8], Jacques S Beckmann[9,10], Sonja I Berndt[11], The Prostate, Lung, Colorectal, and Ovarian (PLCO) Cancer Screening Trial[71], Sven Bergmann[9,12], Amanda J Bennett[3,4], Sheila A Bingham[13], Murielle Bochud[14], Morris Brown[15], Stéphane Cauchi[16], John M Connell[17], Cyrus Cooper[18], George Davey Smith[19], Ian Day[18], Christian Dina[16], Subhajyoti De[20], Emmanouil T Dermitzakis[5], Alex S F Doney[21], Katherine S Elliott[3], Paul Elliott[22,23], David M Evans[6,19], I Sadaf Farooqi[2,24], Philippe Froguel[16,25], Jilur Ghori[5], Christopher J Groves[3,4], Rhian Gwilliam[5], David Hadley[26], Alistair S Hall[27], Andrew T Hattersley[6,7], Johannes Hebebrand[28], Iris M Heid[29,30], KORA[71], Blanca Herrera[24], Anke Hinney[28], Sarah E Hunt[5], Marjo-Riitta Jarvelin[22,23,31], Toby Johnson[9,12,14], Jennifer D M Jolley[5], Fredrik Karpe[4], Andrew Keniry[5], Kay-Tee Khaw[32], Robert N Luben[32], Massimo Mangino[7], Jonathan Marchini[33], Wendy L McArdle[19], Ralph McGinnis[5], David Meyre[16], Patricia B Munroe[5], Andrew D Morris[21], Andrew R Ness[19], Matthew J Neville[4], Alexandra C Nica[5], Ken K Ong[1,2], Stephen O'Rahilly[2,24], Katharine R Owen[3,4], Colin N A Palmer[21], Konstantinos Papadakis[26], Simon Potter[5], Anneli Pouta[31,39], Lu Qi[40], Nurses' Health Study[71], Joshua C Randall[3,4], Nigel W Rayner[3,4], Susan M Ring[19], Manjinder S Sandhu[1,32], André Scherag[41], Matthew A Sims[1,2], Kijoung Song[42], Nicole Soranzo[5], Elizabeth K Speliotes[43,44], Diabetes Genetics Initiative[71], Holly E Syddall[18], Sarah A Teichmann[20], Nicholas J Timpson[3,19], Jonathan H Tobias[45], Manuela Uda[46], The SardiNIA Study[71], Carla I Ganz Vogel[28], Chris Wallace[36], Dawn M Waterworth[42], Michael N Weedon[6,7], The Wellcome Trust Case Control Consortium[72], Cristen J Willer[47], FUSION[71], Vicki L Wraight[2,24], Xin Yuan[42], Eleftheria Zeggini[3], Joel N Hirschhorn[44,48–51], David P Strachan[26], Willem H Ouwehand[5], Mark J Caulfield[36], Nilesh J Samani[33], Timothy M Frayling[6,7], Peter Vollenweider[52], Gerard Waeber[52], Vincent Mooser[42], Panos Deloukas[5], Mark I McCarthy[3,4,73], Nicholas J Wareham[1,2,73] & Inês Barroso[5,73]

## Genome-wide association study identifies eight loci associated with blood pressure

Christopher Newton-Cheh[...] & Patricia B Munroe[...]

---

## Current insights from GWAS:

- Well-powered (meta-)studies with (ten-)thousands of samples have identified a few (dozen) candidate loci with highly significant associations

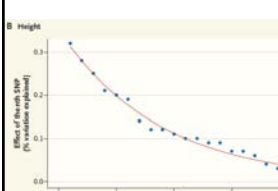- Many of these associations have been replicated in independent studies

## Current insights from GWAS:

- Each locus explains but a tiny (<1%) fraction of the phenotypic variance
- All significant loci together explain only a small (<10%) of the variance

**David Goldstein:**

"~93,000 SNPs would be required to explain 80% of the population variation in height."

*Common Genetic Variation and Human Traits*, NEJM 360;17
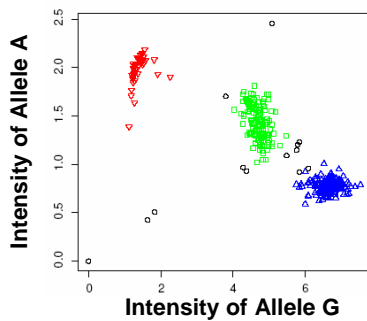
5

## So what do we miss?

1. Other variants like Copy Number Variations or epigenetics may play an important role
2. Interactions between genetic variants (GxG) or with the environment (GxE)
3. Many causal variants may be rare and/or poorly tagged by the measured SNPs
4. Many causal variants may have very small effect sizes
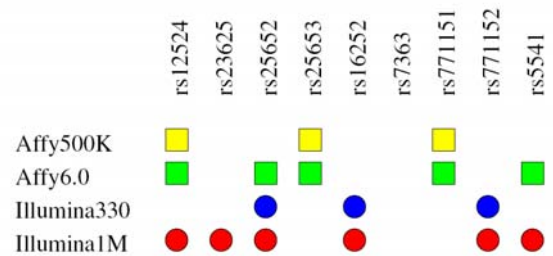5. Overestimation of heritabilities from twin-studies?

## Overview

- Population stratification
- Associations: Basics
- Whole genome associations
- **Genotype imputation**
- Uncertain genotypes
- New Methods
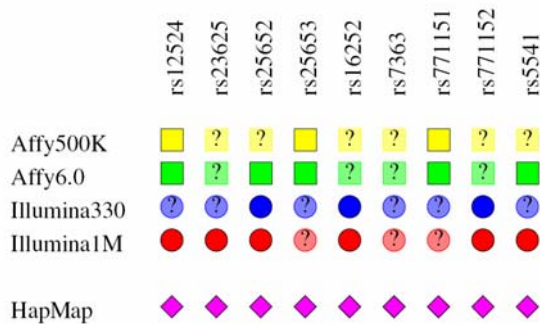
---

### Genotypes are *called* with varying uncertainty



$\nabla$ = AA  $\square$ = AG  $\triangle$ = GG  $\bigcirc$ = not called
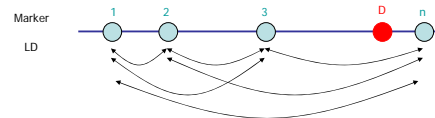
---

### Some Genotypes are missing at all …



---

### … but are *imputed* with different uncertainties



---

### … using Linkage Disequilibrium!



**Markers close together on chromosomes are often transmitted together, yielding a non-zero correlation between the alleles.**

## Conclusion

- Genotypic markers are *always* measured or inferred with *some* degree of uncertainty

- Association methods should take into account this uncertainty

---

## Two easy ways dealing with uncertain genotypes

1. *Genotype Calling:*
   Choose the most likely genotype and continue as if it is true
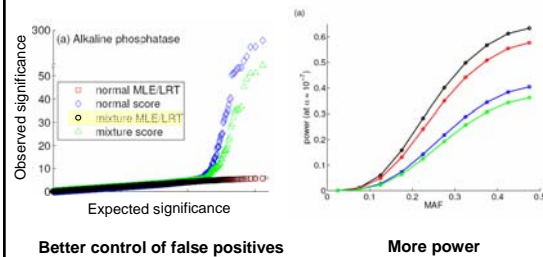   ($p_{11}=10\%$, $p_{12}=20\%$  $p_{22}=70\%$ => $G=2$)

2. *Mean genotype:*
   Use the weighted average genotype
   ($p_{11}=10\%$, $p_{12}=20\%$  $p_{22}=70\%$ => $G=1.6$)
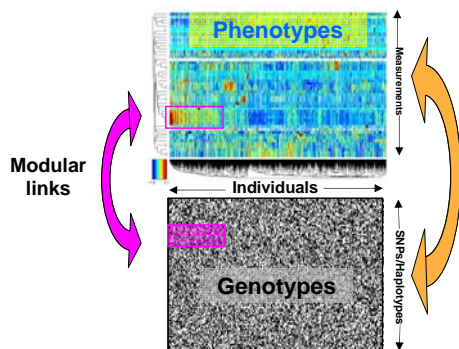
---

## Overview

- Associations: Basics
- Whole genome associations
- Population stratification
- Genotype imputation
- Uncertain genotypes
- New Methods

---

## New Method
based on a mixture model *both* for phenotypes and uncertain genotypes



**Better control of false positives**    **More power**

---

## Modular Approach for Integrative Analysis of Genotypes and Phenotypes



---

## Network Approaches for Integrative Association Analysis



Using knowledge on physical gene-interactions or pathways to prioritize the search for functional interactions

## Overview

- Associations: Basics
- Whole genome associations
- Population stratification
- Genotype imputation
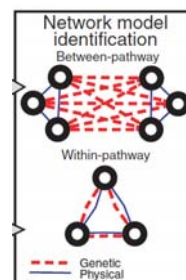- Uncertain genotypes
- New Methods