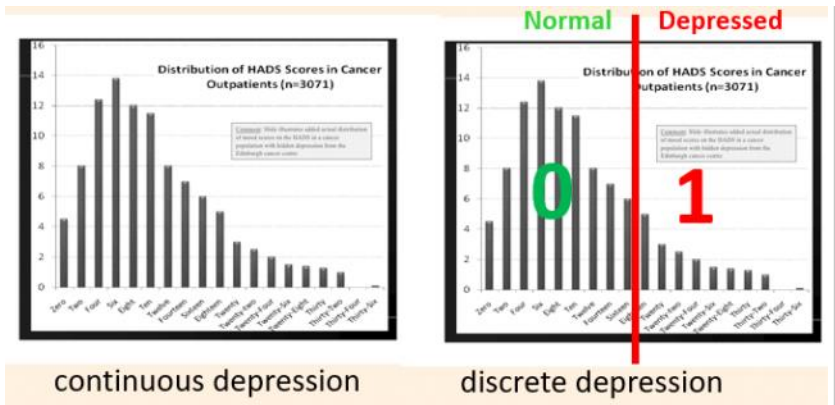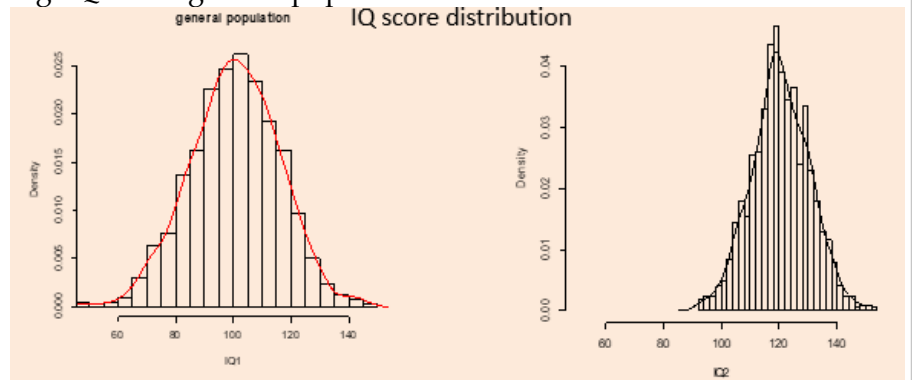# 1a. From GWAS to the twin model via biometrical genetics (c.v.dolan@vu.nl)

| | |
|---|---|
| **What is the commonality between GWAS studies and twin studies?** | Genome wide association studies (GWAS)<br>    The regression of a phenotype on measured variants<br>Twin study<br>    Infering genetic effects from the phenotypic resemblance among twins<br><br>Common biometrical underpinning - Relating genetic differences to phenotypic differences |
| **Why is that 'having five fingers' is not in our genes?** | GWAS is about the analysis of differences<br>    IQ genes - Genes that predicted differences in IQ<br>    Having five fingers is not in our genes - We can compare difference between people with five finger and people with more or less fingers |
| **What is the science of biometrical genetics?** | Biometrical genetics<br>    The science concerned with *inheritance of quantitative traits*<br>    Uses statistical analysis of the inheritance of difference phenotypes as related to plant or animal breeding |
| **Can a continous phenotype (like depression) be discrete?** | Statistics refresher<br>    Linear model - Random variables and parameters<br>    $Y = b0 + b1*x + e$<br><br>    Error is not observed - we do not know what the true values of b0 and b1 are<br><br>    Distribution of random variables - Expressed in a histogram<br><br>    RV are discrete or continuous<br>        ○ Discrete - Dice rolls, number of fingers<br>        ○ Continuous - Height (can be measured with infinitive precision)<br><br>        This varies according to your proposition - Depression could be discrete (you are or you are not) or continuous (there are hundreds of genes related to depression) |

continuous depression    discrete depression

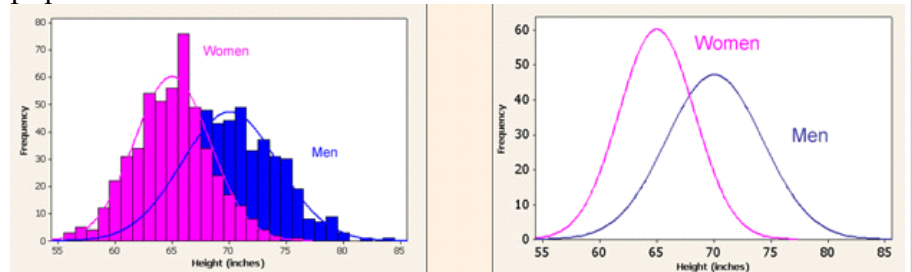| | |
|---|---|
| **Why the distribution of a continuous phenotype is not always normal?** | Distribution of a variable is a property of the different populations<br>    E.g. IQ in the general population and in students<br> |
| **Why is it incorrect to say that 'height is normally distributed'?** | Not all continuous variables are normally distributed<br>    'A variable is normally distributed' -  This is a statement about the population<br><br>    Height is normally distributed is an incorrect statement, since it depends on gender |
| **What is the difference between covariance and correlation?** | Parameters b0 and b1<br>    Mean(y) - b0 + b1*mean(x)<br>    Variance(y) = b1sqrt*variance(x)+ variance e<br><br>Covariance between X and Y - Relationship between two random variables<br>    Pearson Product Moment - Linear association between two continuous variables<br>    $-1 < r < +1$<br><br>Covariance table - Standardized by the correlation coefficient |

| | X | Y |
|---|---|---|
| X | 1 | .350 |
| Y | .350 | 1 |

Z-score - Linear association between X and Y in standardized units

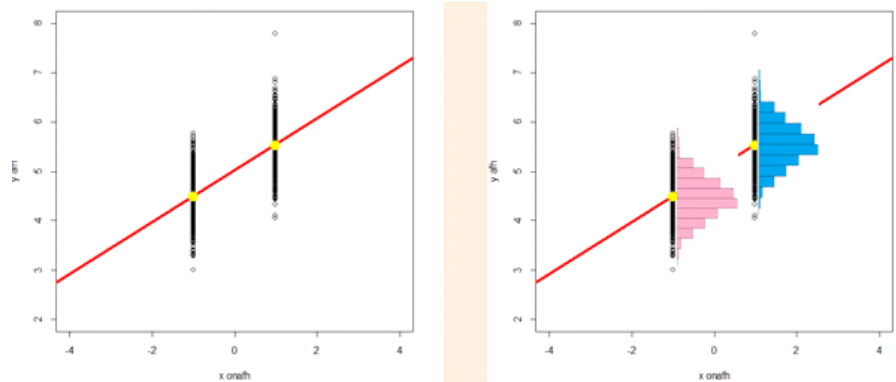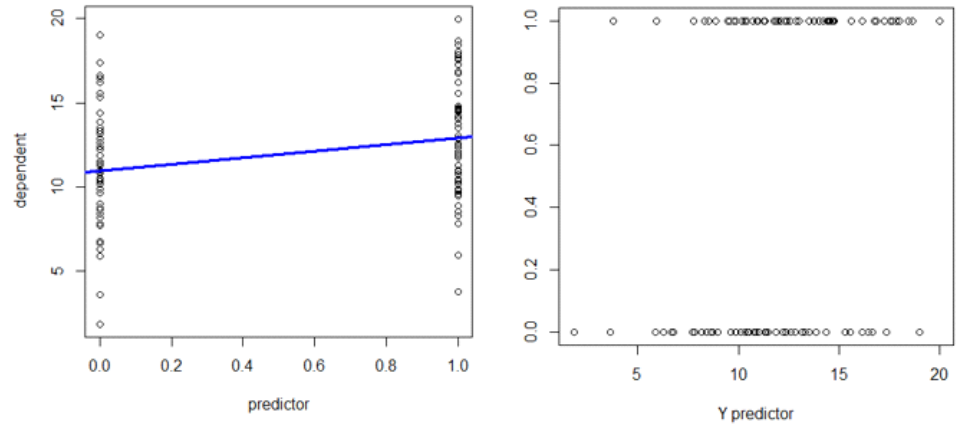| | |
|---|---|
| **What is the difference between regression and correlation?** | Regression - Linear association between X and Y where X is the predictor and Y is the dependent<br><br>    The value of Y is determined by two factors - b1sq*var(x) + var e<br><br>    Are the individual differences of Y related to the individual differences of X?<br><br>    R2 - Amount of variance explained from one variable to another<br>        It is an effect size - It doesn't depend on which variables are used<br><br>    Think about the individual components of the equation |
| **What does it mean to say that a variable is normally distributed given another variable?** | Distributional assumption of statistical inference?<br>    Y I X - Normal (b0 + b1*x, stdev e<br><br><br><br>    Y given X is normally distributed - **In each group of Y, the values of X are normally distributed**<br><br>    It does not mean that Y is normally distributed or that X is normally distributed |
| **What statistical method should we used for** | A binary dependent variable is not suited for linear regression - A logistic regression must be used |

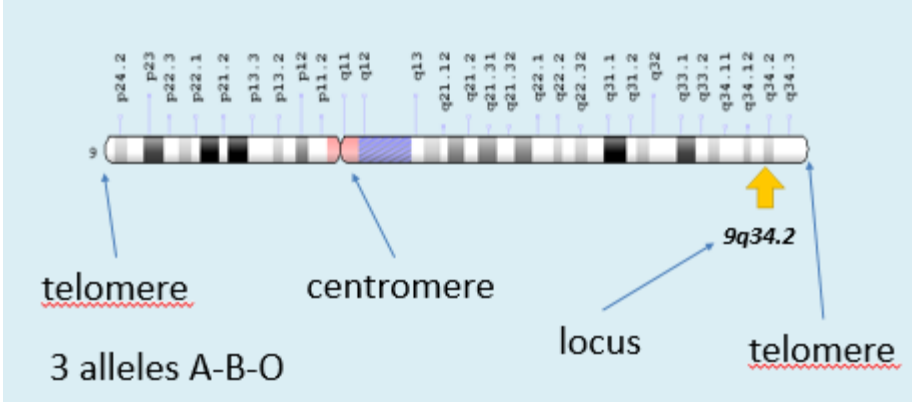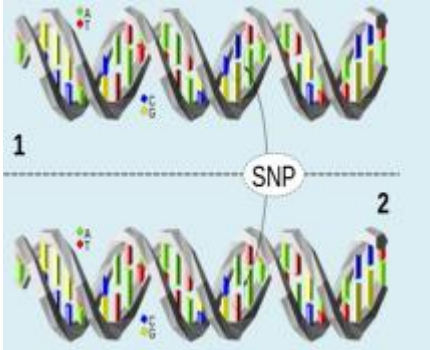| | |
|---|---|
| **binary variables?** | <br><br>This must be analysed with probabilites |
| **What is the difference between conditional probabilities and unconditional probabilities?** | Probabilites - Discrete random variable with a discrete outcome<br>    Unconditional probability<br><br>    Example: Hair and eye color<br>        How many people have light eyes and red hair - Unconditional<br>        How many people have dark eyes - Unconditional<br><br>    Conditional probabilities<br>        What is the probability of light eyes given red hair -<br>        Conditional probability |

# 1b. How do people differ genetically?

| Define: | Terminology |
|---|---|
| a. **Gene** <br> b. **Locus** <br> c. **Allele** | Gene - Sequence of DNA that code for a particular product <br> Locus - Site of a specific gene on a chromosome <br> Allele - Alternative form of a gene at a locus <br> Genotype - The combination of alleles at a particular locus <br> Phenotype - Observed characteristic, trait |
| **Where would a locus be if it was named "9q34.2"** | Chromosome strucutre <br><br>  <br><br> Each locus has 2 alleles - One paternal and one maternal <br>      Mendel's first law |
| **What is a SNP** | SNP (single nucleotide polymorphism) - Variant at a level of a single base pair <br><br>  <br><br>      Variation of the human genome -> difference in protein structure -> phenotypic differences |
| **What are Mendelian traits?** | Mendelian traits - 1 to 1 phenotype-genotype relationship <br>      Variations in 1 gene causes variations in the trait |

Examples
- Sickle cell anemia
- Cystic fibrosis
- Xeroderma pigmentosum
- PKU of fenylketonurie
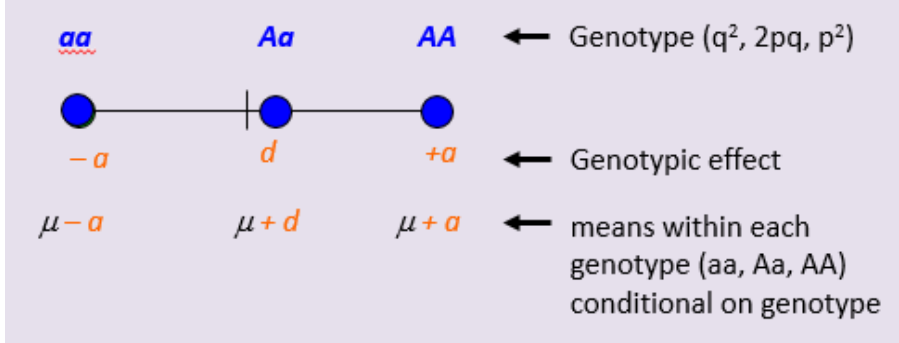- polydactyly

Sickle cell anemia
Cystic fibrosis
Xeroderma pigmentosum
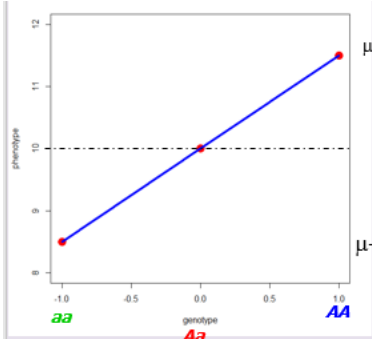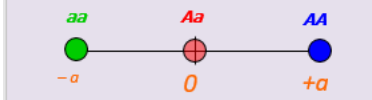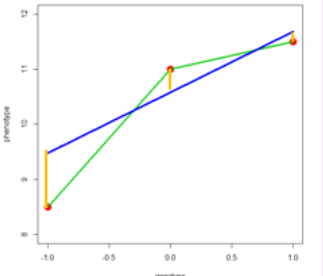Fenylketonurie
Polydactylyl

| What are polygenic traits? | Polygenic traits - Quantitative/complex |
|---|---|
| | Accumulation of many single genes - Quantitative trait loci (QTL) |
| | The accumulation of many different gene give rise to a normal distribution |



| Why is it difficult to have significant results in GWAS? | GWAS cohort study |
|---|---|
| | $Y = b0 + b1*GV + e$ |
| | $Var(y) = B1sqr * var(GV) + var\ e$ |
| | Predictive part depends on B1, the Genetic variance and the error |

| | |
|---|---|
| | H-null: b1 = 0<br><br>Bonferroni correction - With a million tests you are bound to find many false positives |
| **What is Hardy Weinberg equilibrium?** | Define what is the distribution of alleles in the population<br>Biallelic: A and a<br>In GWAS: SNPs<br><br>Frequency of A is p<br>Frequency of a is q<br><br>What are the genotype frequencies (predicted from the the allele frequency)?<br><br><br><br>Hardy Weinberg equilibrium |
| **Why is HW useful for GWAS?** | Observed genotype frequencies - Observed empirically<br>$\quad$ P(AA) = N(AA) / Ntotal<br>$\quad$ P(Aa) = N(Aa) / Ntotal<br>$\quad$ P(aa) = N(aa) / Ntotal<br><br>Estimates allele frequency:<br>$\quad$ P = p(AA) + 1/2 *p(Aa)<br>$\quad$ P = p(aa) + 1/2*p(Aa)<br><br>GWAS studies - Need to establish that what they observe is close to Hardy Weinberg equilibrium<br>$\quad$ This is important because it the genotype testing is not 100% reliable - Some loci are easy to genotype, others are difficult<br>$\quad$ HW assumes that all genotypes have the same fit/ability to reproduce<br>T(1) - Chi-square with one degree of freedom distribution ($p = 3.84$ is |

| | |
|---|---|
| | associated with alpha <0.05) $$T = \frac{(N_O(AA)-N_E(AA))^2}{N_E(AA)} + \frac{(N_O(Aa)-N_E(Aa))^2}{N_E(Aa)} + \frac{(N_O(aa)-N_E(aa))^2}{N_E(aa)}$$ Degrees of freedom - 2x2 pq table; once you know p, you know q |
| **How is the "mean phenotype" measured in a population?** | We know frequencies, we now assign effects to the genotypes M - a - Effect of genotype aa M + d - Effect of genotype Aa M + a - Effect of genotype AA  Take all individuals with a genotype and calculate their mean phenotypes |
| **Describe the formula for the contribution of QTL to the phenotype mean?** | Contribution of the QTL to the phenotype mean m= a(p-q) + 2pqd **to the population phenotypic mean** $\mu$ + m M = a(p-q) + 2pqd A = Homozygous effect D = Heterozygous effect If a and d equal 0 - There is no effect of the genetic variant on the phenotype |
| **What is the difference between additive and dominant effects?** | Additive or linear effects give rise to variance component $s^2_{QTL(A)} = 2*pq[a+(q-p)d]^2$ Dominance or within local allelic interaction effects give rise to variance component $s^2_{QTL(D)} = (2pqd)^2$ Additive effects - d equals 0 |

| | |
|---|---|
| | Dominance effects - d is different from zero |
| **How to derive allele distribution from phenotype?** | Deriving the distribution of alleles in a population:<br>    Take all aa individuals and calculate their mean phenotypic value:<br>    mu - a (conditional mean is an arbitrary variable) |
| | |
| **In a graph of genotype plotted against genotype, what is the visual difference between additive effects and dominant effects?** | When d equals zero, the relationship is dominant - Linear model<br>When d is not equal to zero, the relationship is not dominant - Non-linear model<br><br>Additive/Linear effects<br><br>    Explained variance: ssqr = 2pq(a)sqr<br>        If it equals zero - There is no dominance or recessiveness<br><br>Dominance effects - d is in the middle<br><br>    Explained variance: ssqr = 2pq(a+(q-p)d)sqr<br>        Not explained: ssqr = (2pqd)sqr<br><br>    Dominance becomes part of the residual |

Additive/Linear effects graph labels:

$\mu + a$

Explained variance:

$$s^2_{QTL(A)} = 2pq[a]^2$$

$$s^2_{QTL(D)} = 0$$

$\mu - a$

aa   Aa   AA

$-a$   $0$   $+a$

Dominance effects - d is in the middle graph labels:

$\mu + a$
$\mu + d$

Explained variance (blue line):

$$s^2_{QTL(A)} = 2pq[a+(q-p)d]^2$$

Not explained:

$\mu - a$   $s^2_{QTL(D)} = (2pqd)^2$

aa   Aa   AA

$-a$   $d$   $+a$

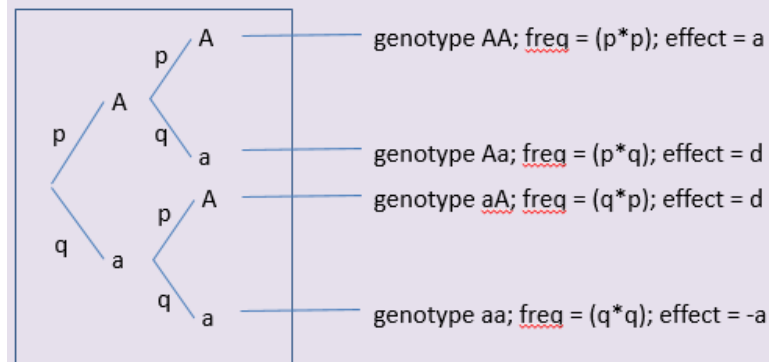| | |
|---|---|
| **What is alpha (in the context of GWAS)?**<br><br>**How to calculate alpha?** | Alpha - Average effect of allele substitution (same as b1)<br>    Derivation:<br>    Image of population (aA and Aa is different)<br><br><br><br>Take all people with large A in the first position<br><br><br><br>       What is the probability that the second allele is A? P (effect a)<br>       What is the probability that the second allele is a? Q (effect d)<br>       Conditional mean1 = p*a + q*d<br>    Take all people with small a in the first position<br>       What is the probability that the second allele is A? P (effect d)<br>       What is the probability that the second allele is a? Q (effect -a)<br>       Conditional mean2 = p*d + q*(-a)<br><br>Alpha = Conditional mean1 - Conditional mean2 |

genotype AA; freq = (p*p); effect = a

genotype Aa; freq = (p*q); effect = d

genotype aA; freq = (q*p); effect = d

genotype aa; freq = (q*q); effect = -a

difference = average effect of allele substitution

$\alpha = \alpha_1 - \alpha_2 = (p*a + q*d) - (p*d+q*-a) =$

pa +qd -pd +qa =

pa +qa - pd + qd =

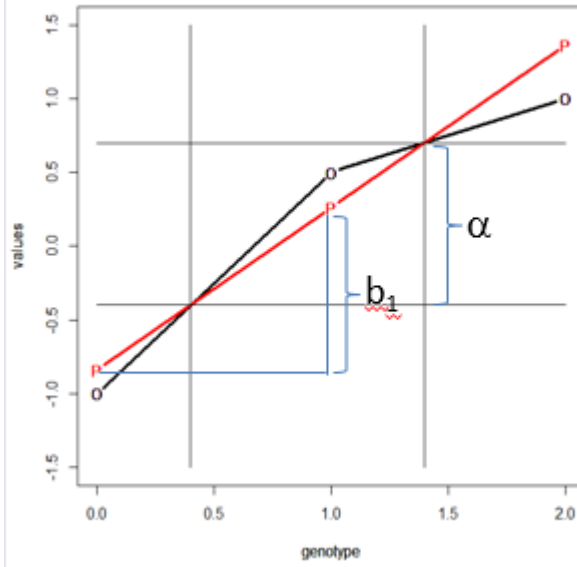$(p+q)a + d(q-p) = a + d(q-p)$

$b_1 = \alpha = (a + d(q-p))$

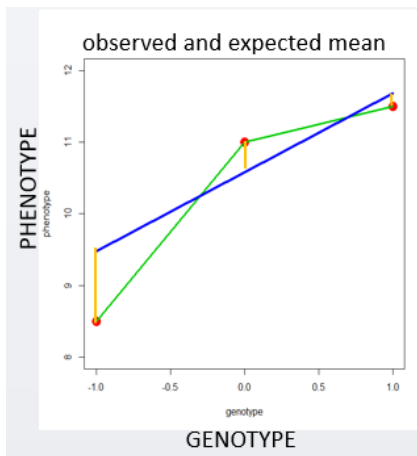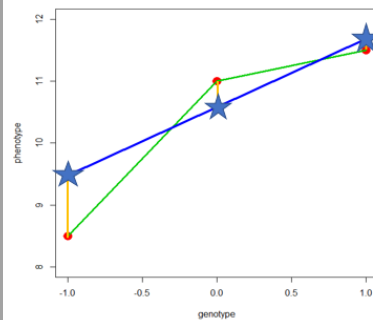| | |
|---|---|
| **Why is alpha the same as b1 in an additive model?** | Predicitive values $$b_1 = \alpha = (a + d(q-p))$$  Genotype AA = 2qalpha<br>Genotype Aa = (q-p)alpha<br>Genotype aa = 2palpha |

# 2. Additive + Dominance Component

| | |
|---|---|
| **How to observe additive or dominant effects in a graph?** | observed and expected mean<br><br>PHENOTYPE / phenotype (y-axis)<br>GENOTYPE / genotype (x-axis)<br><br>Dots - Conditional means (calculated from phenotypic data)<br>Blue line - Regression model<br>Yellow line - Residuals |
| **What are breeding values?** | Breeding values = predicted values from the regression model<br><br>phenotype (y-axis), genotype (x-axis)<br><br>breeding values = predicted values (in the linear regression model) = conditional means, mean(phenotype \| GV)<br><br>Why is it called 'breeding values'? They are predictive of the offspring phenotype and therefore can be used in animal and plant breeding programs |
| **What is the main variable that predicts the success of a breeding program?** | Example: Breeding sheeps<br><br>Goal: Extend fertility period<br>Identify genes that have an effect on extended fertility<br>Choose sheep with highest breeding values<br>The higher the effect of GV on phenotype (higher r squared), the more sucessful the breeding program is |

| | |
|---|---|
| **What is identity by descent? How many alleles IBD do monozygotic twins have? And dyzygotic twins?** | Identity by descent - When the alleles can be traced from the parents<br><br><br><br>Each parent provides one allele for every allele possible<br>Monozygotic twins = 2 alleles IBD<br><br><br><br>Dizygotic twins = Can have 0, 1 or 2 alleles IBD<br>Non twins = Always 1 IBD (one from each parent) |
| **What is the proportion of additive variance shared by**<br>a. **unrelated people**<br>b. **Parent-offspring**<br>c. **MZ twins** | Proportion of shared variance<br><br>$2pq[a+(q-p)d]^2$  $(2pqd)^2$<br>$s^2_{QTL(A)} = b_1^2 s^2_{GV}$  $s^2_{QTL(D)}$<br><br>IBD=0  0  0  Unrelated<br>IBD=1  ½  0  Parent - Offspring<br>IBD=2  1  1  MZ twins<br><br>IBD=0  0  0  25% DZ twins<br>IBD=1  ½  0  50% DZ twins<br>IBD=2  1  1  25% DZ twins |

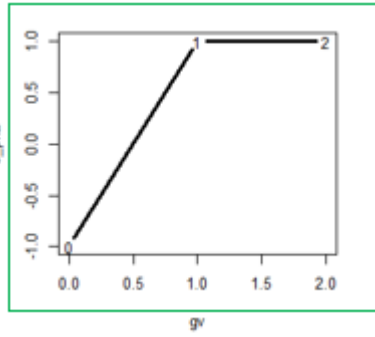| | |
|---|---|
| **What is the proportion of dominance variance shared by**<br>  a.  **unrelated people**<br>  b.  **Parent-offspring**<br>  c.  **MZ twins** | Additive genetic components - Not shared by unrelated people, shared 50% by parent and offspring<br>Dominant genetic components - Not shared by parents and offspring |
| **What is the difference between additive coding, dominant coding and recessive coding?** | Additive coding<br><br>AA = 2<br>Aa = 1<br>aa = 0<br><br>Dominant coding<br><br>AA = 1<br>Aa = 1<br>aa = 0<br><br>Recessive coding |

AA = 1
Aa = 0
aa = 0
Figure - Dominant coding and recessive coding

| **How to represent non-linear relationships in a regression model?** | What about non-linear relationships? |

| co-dominant coding | | |
|---|---|---|
| genotype | $allele_1$ | $allele_2$ |
| $(A_1A_1)$ | 1 | 0 |
| $A_2A_1$ & $A_1A_2$ | 0 | 1 |
| $A_2A_2$ | 0 | 0 |

Dummy coding - Allelic association test

| **How is the additive and dominant effects of genotype described in model to explain variance of phenotype?** | Explore genotype-phenotype relationship |

| Polynomial | | |
|---|---|---|
| | additive | dominance |
| genotype | $GV_A$ | $GV_D$ |
| $A_1A_1$ | 2 | 4p-2 |
| $A_1A_2$ $A_2A_1$ | 1 | 2p |
| $A_2A_2$ | 0 | 0 |

better coding for studying genetic model:
var(pheno) = $b_{1A}^2$*var($GV_A$) + $b_{1D}^2$*var($GV_D$) + var(residual)
var(pheno) = 2pq[a+(q-p)d]²+ (2pqd)²      + var(residual)

Reject null hypothesis - Dominance is present
Coding guarantees that the correlation between Gva and Gvd is zero (avoid multicollinearity in the model)

| **How can dominance be visually identified from a graph?** | Graph - a big difference between the lines indicates dominance |

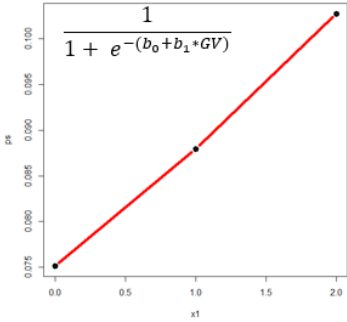| | |
|---|---|
| **What is the difference between a simple linear model and a logistic regression?** | Dichotomous outcome variable - A logistic regression model must be used<br>    Probability of scoring 1 (having the disease) based on the genotype<br>variant<br><br>**logistic function** $prob(1\|GV) = \dfrac{1}{1+ e^{-(b_0+b_1*GV)}}$<br><br>    Exp (x) = e raised to the power of x |
| **What happens when there is a big effect of genetic relative risk in a logistic regression model?** | Genetic relative risk - Ratio of risk based on a reference<br>(condGV=1/condGV=0)<br><br><br><br>1 = no effect<br>Large = big effect (the regression line becomes S-shaped) |

# 3a. Linkage equilibrium and disequilibrium

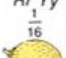| | |
|---|---|
| **What is Mendel's first law of segregation?** | Mendel's first law of segregation - Parents alleles are chosen at random<br>    Random selection - 1 out 4 for each parent<br><br> |
| **What is Mendel's second law of segregation?**<br><br>**When is Mendel's second law of segregation violated?** | Mendel's second law of segregation - The assortment for each allele is independent from one another<br><br><br><br>Exception; Linkage - Proximity of genetic loci in a chromosome |

| | Creates gametic phase disequilibrium |
|---|---|
| **What is recombination?** | Crossing over and recombination<br>   Crossing over happens during second phase of meiosis<br><br><br><br>   Chromosomes A and D - No crossing over<br>      Alleles a and c - Violates Mendel's second law because they are transmitted together<br>   Chromosomes B and C - Only happens due to cross over |
| **What happens when two loci are very close together regarding Mendel's second law? And when they are in different chromosomes or very far apart in the same chromosome?** | Probability of recombination depends of the distance<br>   Large distance - Probability is 0.5 (Mendel's second law holds)<br>   Short distance - Probability is very close to 0<br><br><br><br>ABD is a haplotype<br>ABd is a new haplotype |

| What is theta? | Same chromosomes - Theta is the recombination probability |
|---|---|
| | Theta is a function of distance - Base pairs of DNA |



If theta is 0 - Recombination is unlikely
      Once you observe one allele, you are likely to observe the other
      (markers can be used to observe causal genes)



If theta is 0.5 - Recombination is as likely as non-recombination

| **What is the value of theta when the loci are in different chromosomes?** | Different chromosomes - Theta is 0.5 (chromosomes are chosen at random, there is no dependency) |
|---|---|

| | |
|---|---|
| |  chromosome 1 chromosome 2 a b c d $\theta = 0$ k l m n $\theta = 0$ |
| **What is the unit of theta?** | Distance in chromosome - Expressed in centimorgans 1 Morgan = 1 million basepairs 1 centiMorgan = 10 thousand basepairs  Theta - Depends on centimorgans (distance between loci) |
| **If the haplotypes depend on allele probability, what is the value of theta?** | Quantifying close: Linkage equilibrium Haplotypes depend on the allele probability Loci are independent - Either very far apart in the same chromosome or in different chromosomes (theta = 0.5) |

| | |
|---|---|
| | Alleles of marker 2 ($p_B$ is freq.)<br><br>**B ($p_B$)**  **b ($p_b$)**<br><br>**Alleles of marker 1 (freq.)**<br><br>**A ($p_A$)**<br><br>AB<br>($p_A p_B$) \| Ab<br>($p_A p_b$)<br><br>**a ($p_a$)**<br><br>aB<br>($p_a p_B$) \| ab<br>($p_a p_b$) |
| **What is the influence of linkage disequilibrium and HW disequilibrium?** | Linkage disequilibrium between two markers<br>    Has no bearing on Hardy Weinberg equilibrium (genotypes at each<br>    marker do not have to be randomly selected) |
| **What is D'? What happens when D' equals 0?** | Measure of LD: D′ (standardization value that indicates linkage disequibrilium)<br><br>AB<br>$P_{AB}=(p_A p_B+D)$ \| Ab<br>$P_{Ab}=(p_A p_b-D)$<br><br>aB<br>$P_{aB}=(p_a p_B-D)$ \| ab<br>$P_{ab}=(p_a p_b+D)$<br><br>If D′ is not 0, the loci are correlated (Linkage equilibrium)<br>If D′ is different from 0, loci are not perfectly correlated (Linkage disequilibrium)<br><br>A value of D′ on its own is not very informative |
| **What is the cause of linkage** | Linkage disequilibrium decays over generations as a function of theta |

| | |
|---|---|
| **disequilibrium in nature? And linkage equilibrium?** | Disequilibrium happens due to a mutation<br>Equilibrium happens due to crossing over<br><br>The higher the theta, the faster disequibrilium decays<br>The smaller the theta, the less likely it is to break the loci apart<br><br>If the loci are in two different chromosomes, it only takes one generation to reach equilibrium |
| **What is r squared related to Linkage Disequilibrium? How can it be calculated?** | R squared - Direct function of LD<br> Same as correlation between two loci squared<br><br>R square is 0, LE<br>R square is 1, LD |
| **Why should we use markers in GWAS studies instead of measuring the disease locus directly?** | Direct and indirect allelic association<br> D = Disease locus<br><br>Direct - Measure disease relevance directly<br> Usually not possible because you don't know what it is<br><br>Indirect - Measure markers to assess trait effects on D |

The graph (image 1) shows curves labeled $\theta = 0.0001$, $\theta = 0.001$, $\theta = 0.01$, $\theta = 0.1$, $\theta = 0.5$, with y-axis from "Complete disequilibrium" to "Complete equilibrium" and x-axis "Generations" from 1 to 1000.

© Elsevier. Nussbaum et al: Thompson and Thompson's Genetics in Medicine 7e - www.studentconsult.com

| | |
|---|---|
| | Hits on a GWAS does not mean much -> you still don't know what D is, only a range of where it can be -> posterior studies are necessary<br>Significance of markers (red squares) - Indicate upper and lower boundaries of where the disease locus is<br> |
| **How can GWAS studies be significant if the effect of each QTL is small and the alpha level is tiny?** | GWAS - Measure GV and phenotype -> linear or logistic analysis<br>    Equation<br><br>    $E[pheno\|GV] = b_0 + b_1 * GV$    (conditional mean or predicted value)<br>    $E[pheno=1\|GV] = 1 / \{1\text{-}exp(-(b_0 + b_1 * GV))\}$  (conditional probability)<br><br>If effect size is small and alpha is small, you need a very large sample size to pick up values |
| **What happens when you included marker and disease locus in a regression?** | Example<br>    Red lines - Not significant<br>    Green lines - Significant<br>    S - Markers<br>    Q - Causal loci to phenotype<br><br>If Q is included, markers will not explain anything (Q absorbs)<br>If Q is not included, markers will explain something (because they are correlated with Q) |

| | |
|---|---|
| **What are the three common forms of association in GWAS?** | Association in GWAS: 3 common forms<br>    Direct association - Measure of SNP of interest that is causing the disease<br>    Indirect association - Markers<br>    Spurious association - You reject the null hypothesis, but the SNP is not related at all |
| **What is a probable cause of spurious associations in GWAS?** | Cause of spurious associations: Population stratification<br>    If you put two different populations together, you may get significant results |

# 3b. Classical twin design

| | |
|---|---|
| **Why are twin studies used to study genetic effects?** | Example: CHRM2 gene - Feedback and regulation of Ach release, related with higher cognitive processing<br>      Find a candidate gene first<br>      Posterior study confirmed the results |
| **Why are GWAS meta-analyses becoming more and more common?** | GWAS meta-analysis<br>      Compendium of dozens of studies - Increase sample size<br><br>GWAS opens the black box - One SNP at a time<br>      There are too many QTLs for complex traits like IQ<br>      Variance of IQ = additive effects + dominant effects + error<br>      (environmental effects) |
| **What is shared variance?** | Shared variance - If two variables are subject to the same influence, they share variance attributable to that influence<br><br><br>      Covariance - Twins share genetic variance |
| **Why are dyzygotic twins share 25% of their genetic variance IDB?** | Genetic resemblance - Not affected by recombination/crossing over<br>      MZ twins are genetic identical - IBD 2<br>      DZ twins - IBD 0, 1 or 2 -> average of 1/4<br>      Half siblings<br> |

| Table 2.18. (1.24): shared genetic variance | | |
|---|---|---|
| Relationship | Genetic variance component | |
| | $s_A^2$ | $s_D^2$ |
| DZ twins, full sibs* | ½ | ¼ |
| MZ twins | 1 | 1 |
| Half-siblings* | ¼ | 0 |
| Parent-Offspring | ½ | 0 |
| Unrelateds | 0 | 0 |
| informative w.r.t. variance sharing | proportion of alleles shared IBD $s_A^2$ sharing | probability of IBD=2 sharing $s_D^2$ sharing |

| | |
|---|---|
| **Which answer does the classical twin design try to answer?** | Classical twin design<br>    How many variance of additive and dominance effects depends on genetic relatedness<br>    Explained variance of the genetic variant with relation to the phenotype = coefficient component + variance of the genetic variant |
| **Why is pathdiagram notation useful?** | Pathdiagram notation - Regression model<br>    Square - Observed variable<br>    Circle - Non-observable variable<br>    Double headed arrow - Variance of X<br><br>latent variable Y (not observed)<br>latent variable X (not observed)<br>$S_Y^2$ variance of Y<br>observed variable Y<br>$S_X^2$ variance of X<br>latent variable X )<br>Images<br>Covariance<br>$S_X^2$   $b_1$   X → Y ← 1 ε   $S_ε^2$ |

| | |
|---|---|
| | Regression (direction of effect - single arrow)<br><br><br><br>Regression with two terms - Term that represents correlation between the variables (2.b1.b2) |
| **Draw a pathdiagram of MZ twins.** | Pathdiagram of monozygotic twins<br><br><br><br>Seconds representation is clearer<br>D - Hypothetical latent variable<br>IBD = 1 |
| **Draw a pathdiagram of dyzygotic twins.** | Pathdiagram of dizygotic twins |

For the MZ twins diagram, the labels include: $1*\sigma_A^2$, $\sigma_D^2$, $\sigma_D^2$, $\sigma_A^2$, A, D, D, A, $\sigma_A^2$, 1, 1, $1*\sigma_D^2$, 1, 1, phenotype MZ$_1$, phenotype MZ$_2$, $\varepsilon_1$, $\varepsilon_2$.

| | MZ twin 1 | MZ twin 2 |
|---|---|---|
| **MZ twin 1** | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | $1*\sigma_A^2 + 1*\sigma_D^2$ |
| **MZ twin 2** | $1*\sigma_A^2 + 1*\sigma_D^2$ | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ |

beh gen – Lecture 3  2018

|  | DZ twin 1 | DZ twin 2 |
|---|---|---|
| DZ twin 1 | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | $\frac{1}{2}*\sigma_A^2 + \frac{1}{4}*\sigma_D^2$ |
| DZ twin 2 | $\frac{1}{2}*\sigma_A^2 + \frac{1}{4}*\sigma_D^2$ | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ |

IBD = 1/2

| **How can you calculate the covariance matrix of twins?** | Covariance matrix of twins |
|---|---|

|  | MZ twin 1 | MZ twin 2 | observed | |
|---|---|---|---|---|
| MZ twin 1 | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | $\sigma_A^2 + \sigma_D^2$ | 135.61 | 117.10 |
| MZ twin 2 | $\sigma_A^2 + \sigma_D^2$ | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | 117.10 | 144.12 |

|  | DZ twin 1 | DZ twin 2 | observed | |
|---|---|---|---|---|
| DZ twin 1 | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | $\frac{1}{2}*\sigma_A^2 + \frac{1}{4}*\sigma_D^2$ | 144.99 | 52.79 |
| DZ twin 2 | $\frac{1}{2}*\sigma_A^2 + \frac{1}{4}*\sigma_D^2$ | $\sigma_A^2 + \sigma_D^2 + \sigma_s^2$ | 52.79 | 133.37 |

| $\sigma_A^2$ | $\sigma_D^2$ | $\sigma_s^2$ | $\sigma_{Ph}^2$ |
|---|---|---|---|
| 93.766 | 22.814 | 22.753 | 139.333 |

| expected MZ | | expected DZ | |
|---|---|---|---|
| 139.333 | 116.580 | 139.333 | 52.587 |
| 116.580 | 139.333 | 52.587 | 139.333 |

Monozygotic twin
Covariance = QTLa + QTLd

Dizygotic twin
Covariance = 1/2QTLa + 1/4QTLd

Expected model should resemble the observed effects

| | |
|---|---|
| **What does A, D, C and E stand for in a twin study model?** | A and D are random variables<br><br><br><br>$$\sigma_{Ph}^2 = \sigma_A^2 + \sigma_D^2 + \sigma_s^2 \qquad \sigma_{Ph}^2 = \sigma_A^2 + \sigma_D^2 + \sigma_c^2 + \sigma_E^2$$<br><br>A - Additive effects<br>D - Dominance effects<br>E - Environmental variation<br>    C - Shared environmental variants<br>    E - Non-shared environmental variants<br><br>Error is not included in the model |
| **How are addivite, dominant and environmental effects calculated in a twin study?** | Additive + Dominant = Genetic component (broad sense heritability)<br>    Overall effect<br>Additive only - Narrow sense heritability<br><br>$\sigma_A^2 \qquad\qquad \sigma_D^2 \qquad\qquad \sigma_s^2$<br>$4*r_{DZ}-r_{MZ} \qquad 2*r_{MZ}-4*r_{DZ} \qquad 1- \sigma_A^2 - \sigma_D^2$<br>$4*.18 - .41 \qquad 2*.41-4*.18 \qquad 1-.31-.10$<br>$.31 \qquad\qquad\quad .10 \qquad\qquad\quad .59$<br><br>QTLa = 4*rdz-rmz<br>QTLd = 2*rmz-4rdz<br>QTLe = 1 - QTLa - QTLd (i.e. the rest) |
| **What is the main problem with the classical twin design?** | Twin design issue<br>    Distiction between genetic and environmental effect - You cannot fit a model with additive, dominant, shared environmental effects and unshared environmental effects<br>    ACDE model is not possible |

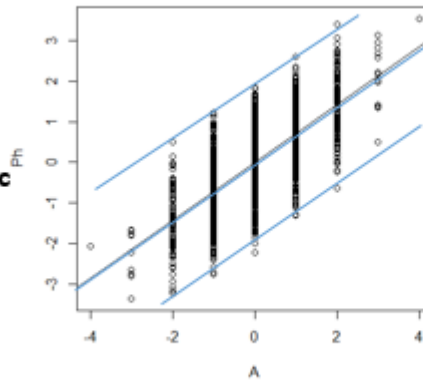| | |
|---|---|
| | ACE model is ok<br>ADE model is ok<br><br>E model = twins are not correlated<br>CE model = monozygotic and dyzygotic twins are the same |
| **What are the four types of models possible in a twin study?** | ADE model – estimate variance components from twin correlations<br><br>| | $\sigma_A^2$ | $\sigma_D^2$ | $\sigma_s^2$ |<br>|---|---|---|---|<br>| calculation | $4*r_{DZ}-r_{MZ}$ | $2*r_{MZ}-4*r_{DZ}$ | $1-\sigma_A^2-\sigma_D^2$ |<br><br>ACE model – estimate variance components from twin correlations<br><br>| | $\sigma_A^2$ | $\sigma_C^2$ | $\sigma_E^2$ |<br>|---|---|---|---|<br>| calculation | $2*(r_{MZ}-r_{DZ})$ | $2*r_{DZ}-r_{MZ}$ | $1-\sigma_A^2-\sigma_D^2$ |<br><br>If rMZ > 2*rDZ = ADE<br>If rMZ < 2*rDZ = ACE<br>If rMZ = 2*rDZ = AE model<br>If rMZ = rDZ = E model<br><br>When the difference is not very large -> Small dominance component |
| **Why twin studies challenged the status quo in the 80s?** | Twin studies - Highlight the importance of genetic variation<br>    In the 80's, people denied the influence of genetics for complex behaviors |
| **What are the assumptions of the classical twin design?** | Assumptions of CTD<br>• The sample is representative of the population<br>• Random mating - No correlation between the parents, which is often not true (high IQ children have high IQ parents)<br>• Equal environmental assumption<br>• No correlation between A-C or A-E - Safe assumption for personality, risky assumption for IQ (IQ correlates with high economic status, affects the offspring environment)<br>• No Genetic-environment correlation - Younger children intelligence are more dependent on shared environmental effects, older children are more dependent on genes |
| **What is the difference when there is AxE** | No AxE correlation |

| | |
|---|---|
| **correlation and when there is not?** | ɪraction <br><br>AxE correlation <br><br>Depending on A score, environmental effects change |
| **What can you do if the premises of the classical twin design are violated?** | What you can do regarding violation of the classic twin design:<br>• Extend models (include parents)<br>• Extend assumption - Include more information<br>• Violations of assumption are well known -> Data can be interpreted in a specific light |

# 4. From heritability to functional experiments

| Why is the field of human genetics developing so rapidly? | Field of human genetics is rapidly advancing in the recent past<br>    New technologies<br>    Large scale collaboration<br>    Novel disease insights |
|---|---|
| What was interesting to notice regarding the number of publications of Classic Twin Studies in the last few decades? | Meta-analysis of all twin studies<br>    Between 1900-2012<br>    Extraction of sample size, effect size<br>    Standardize trait quantification -ICF<br><br>There is still an increase of twin studies every year, even though we now can genotype individual very cheaply<br><br>Nstudies<br>2,748 studies published between 1958-2012<br><br>Nstudied traits<br>Reporting on 17,804 traits<br><br>Avg. NPairs<br>14,558,903 partly dependent twin pairs<br><br>Most studied traits: Weight, diseases, big-five personality traits, intelligence |
| How much of 'human behavior' is genetically inherited? | If you lump all traits together into 'human behavior'<br><br>Rmz = 0.636<br>RDZ = 0.339<br>Genetic inheritability = 0.49 |

| | |
|---|---|
| | Shared environmental effect = 0.17 |
| **What were the main conclusion of meta-analysis of twin studies?** | Main conclusion<br>    All traits are heritable to some extent<br>    Shared environmental effect is relatively small<br>    Most traits seem to have additive genetic effects<br><br>    Interactive website: match.ctglab.nl |
| **What is heritability?** | Heritability - Proportion of trait variance attributable to genetic variance<br>    Variation in genes underlie trait differences between individuals |
| **How close are humans to:**<br>  a. **Mice**<br>  b. **Chimpanzees**<br>  c. **Humans of the opposite gender** | If heritable - Look at the DNA<br>    3 billion base pairs - One in every thousand base pair is different between people<br><br>    Human and mice - 87,5%<br><br><br><br>Human and chimpanzee - 99%<br><br><br><br>Males and females - 99,9%<br><br><br><br>*50% in twins = Means 50% of 0.1%* |
| **What is the difference between monogenic and** | Monogenic disorders - One gene/one mutation<br>    Most genetic causes are already known<br><br>Polygenic disorders - Influenced by multiple genes, each of small effect |

| | |
|---|---|
| **polygenic disorders?** | Genetic causes mostly unknown (correlation with environmental factors)<br>Genetics variations (SNPs): Inside or outside genes<br>Genetic variations can be: harmless, harmful, latent, silent |
| **How were genetic studies made before the advent of GWAS?** | Discovering genes for complex traits<br>    Candidate gene study - Before 2005<br>    GWAS - Genotyping became cheap and accessible |
| **What was the first GWAS hit discovered?** | GWAS results<br>    First hit discovered - Macular degeneration (a few genes with large effect sizes)<br>    Most cases - Small effect sizes, large samples are needed<br>        Most genetic variance is not in the genes! Gene deserts or intronic regions |
| **What do GWAS hits signify for researchers?** | GWAS hits are hard to interpret<br>    SNPs are correlated in the same chromosome<br>    There are often hundreds of hits - it is difficult to choose a gene for a follow up study |
| **What is a Manhattan plot?** | Manhattan plot - GWAS results<br> |
| **Why molecular studies can help us understand GWAS hits?** | Interpreting GWAS loci<br>    Interference with DNA folding<br>    Positional information - SNP can indicate genes of interest |

| | |
|---|---|
| **What is scRNAseq?** | Novel techniques:<br>scRNAseq - Gene expression at a cellular resolution (GWAS for specific cell types)<br> |
| **How to design follow up experiments from GWAS hits?** | Designing follow up experiments<br>Classical techniques (knockout, knockdown) - not viable for SNPS with small effect sizes<br>Taking into account polygenic nature of traits: Chemo and ontogenetic techniques; induced pluripotent stem cells from patients with disease (create mini-brains in the lab) |
| **What is the monozygotic and dizygotic twin correlation for IQ? Which model does this suggest?** | Individual differences in IQ<br>Twin correlation for IQ - 0.82 MZ and 0.45 DZ |

GWAS - 81 missense mutations associated with IQ
Mostly expressed in the brain
Expressed more in certain cell types - Medium spiny neurons,
pyramidal CA1 neurons
Correlation with different traits

*Score of GWAS cannot predict an individual's IQ score, only to distinguish individuals at a population level*

# 5. GWAS and interpretation of GWAS results

| | |
|---|---|
| **What is the purpose of Disease Genetics?** | Disease Genetics<br>Predict if someone will get sick<br>Test hypothesis about relationships to other diseases or traits (comorbidity)<br>Understand the biology of the disease so we can design better treatments and diagnostics |
| **What is the main difference between GWAS and twin studies?** | Twin and family studies do not measure the variation at all - They inform the contribution of genetics/environment |
| **Why some GVs are correlated to one another?** | What creates genetic variation<br>Mutation - Spontaneous change in the DNA<br>Recombination - Re-shuffles existing patterns of variation<br><br>Consequences - Genetic variants are correlated because they share a history of inheritance<br>Correlation decreases over generations due to recombination<br>Depends on linkage disequilibrium |
| **What did the HAPMAP project do? Why was this important?** | Linkage disequilibrium map of the human genome<br>HAPMAP projects - Genotyping of different populations with different ethnic backgrouns<br>Limited haplotype diversity - Correlated blocks of SNPs<br><br>Only three options for a block of 6 SNPs |

| | |
|---|---|
| **Why was the HAPMAP project important to increase genotype efficiency?** | Genotyping a single SNP in a block provides all the information for all SNPs in this block - Efficient genotyping, less expensive  *Genotyping is always done in the plus strand* |
| **What is the purpose of reimputing SNPs in a GWAS? What is the problem with this process?** | Different companies genotype different SNPs - may be problematic when performing a meta-analysis  You can reimpute the missing correlated SNPs<br>This misses rare haplotypes or recent mutations in the population |
| **What are the steps in GWAS?** | Steps in GWAS until Manhattan plot<br>    Quality of genotyping is important - Heterozygous can be similar to homozygous<br>        Tests for genotyping error<br>    Imputation<br>    Association analysis - PLINK software<br>    Meta-analysis - Increase sample size<br>    Manhattan plot |

| | |
|---|---|
| **What is the common GWAS threshold?** | Statistical tests for every SNP -> High number of false positives<br>    GWAS threshold is 5*10e-8<br>    GOLD standard for association studies - Replicating association in different laboratories |
| **What are examples of false replications and the arguments that were used to try to justify them?** | Not true replications - and proposed explanations<br>• Association to the same trait, but a different gene - Justified by Genetic heterogeneity (different populations)<br>• Association to same trait, same gene, different SNPs - Justified by Allelic heterogeneity<br>• Association to same trait, same gene, same SNP, opposite direction - Justified by Allelic heterogeneity/population differences<br>• Association to different, but a different correlation phenotype - Justified by Phenotypic heterogeneity (depends on phenotype correlation!)<br>• No association at all - Justified by "Sample size too small" |
| **What is a true replication of a GWAS?** | True replication - Same trait, same SNP, same allele, same direction of effect, different and independent population |
| **What do single outlier dots in a GWAS might indicate?** | Manhattan plot<br><br>'Line' is more reliable than single dots - Probably caused by genotyping error |
| **How is causation proved from a GWAS analysis?** | Functional studies - Where is this gene being expressed?<br>    Develop ideas for posterior causal studies<br>    Causation is proved with clinical trials/pharmacological assays |
| **Why is it difficult to interpret GWAS hits?** | Post GWAS annotation and interpretation<br>    LD complicates the interpretation of results - You do not measure causal genes<br>    ○ Some SNPs are known to have no effect<br>    ○ Some SNPs are known to have a direct effect - Change RNA structure or protein; can be located in exons or introns |

| | |
|---|---|
| **What are the functional categories of SNPs?** | Functional categories of SNPs<br>    Protein coding - May alter protein structure (truncated proteins)<br>    Splicing regulation - Disrupt splicing regultation<br>    Transcriptional regulation - Disrupt gene regulation (TF binding<br>    sites, CpG islands, microRNAs)<br>    Post-translation modification - Interfer with proper posttranslation<br>    modification of proteins |
| **If there are many GWAS hits, how can you choose SNPs for posterior studies?** | How to pinpoint causal genes based on GWAS risk loci<br>    Are there functional variants in GWAS risk loci?<br>    Are there SNPs with high CADD scores or low regulomeDB scores?<br>    Are there regulatory variants or eQTLs in GWAS risk loci?<br>Example<br>    CADD scores higher than 10 - Select important SNPs<br>    RegulomeDB score 1f<br>    HiC interaction in the Brain<br><br> |
| **What are CADD scores?** | CADD scores - Combined Annotation Dependent Depletion score<br>    Tool for scoring deleteriousness of SNPs + insertions and deletions in<br>    the human genome<br>    Higher than 10 - Predicted to be 10% most deleterious substitutions<br>    Higher than 20 - Indicates the 1% most deleterious |
| **What are eQTLs?** | Expression QTLs (eQTLs) - Increase level of expression of other genes |

| | |
|---|---|
| | The same regulatory regions and variant could be an eQTL for gene 2 in (a) tissue 1 and for gene 1 in (b) tissue 2, suggesting that limited interrogation of tissues would be misleading for the biological signal underlying disease. <br><br> (a) <br><br>  <br><br> Can be present outside/farther away from the gene of interest |
| **Why studying chromatin interaction may be useful for GWAS?** | Chromatin interaction <br><br>  <br><br> DNA sequences may be distant in the linear sequence but close in the tridimensional structure |
| **What can you study after you find a hit on GWAS?** | SNP annotation implicates genes - Guide for posterior studies <br>      Explore gene functions <br>      Explore pathways of implicated genes <br>      Explore in which tissues the genes are expressed <br>      Explore which cell types are indicated |
| **What is FUMA?** | FUMA (Functional Mapping and Annotation) - Tool for GWAS annotation <br>      Q-Q plot - Red line equals no association (calculated from your alpha level) <br><br>  <br><br>       If every SNP is associated - Suggests population stratification <br> eQTL - Independent genes (linkage equilibrium) |

# 6a. Gene-set analysis

| | |
|---|---|
| **What are the options to test for functional clustering of SNPs?** | Testing for functional clustering of SNP associations  Single SNP analysis<br>Gene-based analysis - Gene as an unit of analysis<br>Gene-set analysis (binary) - Biological pathway<br>Gene-property analysis (continuous) - Expression levels or probability of being a member of a gene-set |
| **What is the problem with testing for the joint association of all SNPs? How would you solve that?** | Gene based analysis<br>• Tests for joint association effects of all SNPs in a gene, taking into account LD (correlation between SNPs -> results in false positives ->multicollinearity must be avoided) (**Exam question**)<br>No single SNP needs to reach GWAS significance -> the gene-based may still be significant |
| **What does a single dot in a gene Manhattan plot might represent?** | SNP Manhattan Plot<br>Every dot is a single SNP  Gene Manhattan Plot<br>Every dot is a single gene  Single dot - Implies that the individual SNPs have a really small effect, but their joint effect has a large effect |

| | |
|---|---|
| **What are the pros and cons of gene-based analysis?** | Gene-based analysis<br>    Pros: Reduce multiple testing, accounts for heterogeneity in gene (small effects of single SNPs), immediate gene-level interpretation (easier to come up with follow-up hypothesis)<br>    Cons: Disregards regulatory/non-genic information, still a lot of tests |
| **What are the pros and cons of using genes as a unit of analysis?** | Using gene as an unit of analysis<br>Pros<br>    Reduce multiple testing<br>    Increases statistical power<br>    Deals with genic heterogeneity<br>    Provides immediate biological insights<br>Cons<br>    Selecting reliable sets of genes is difficult<br>  • Different levels of information<br>  • Different quality of information |
| **On what basis can you choose gene sets?** | Choosing gene-sets can be based on:<br><br>    Protein-protein interaction networks<br>        Yeast-two hybrid or immunoprecipitations - Real differentiated cells, extract real interactions!<br>    Co-expression - Protein interactions expressed at the same in the same tissue/organelle<br>    Transcription regulatory networks<br>    Biological pathway |
| **What is the problem with public databases of gene sets?** | Choosing gene sets<br>    KEGG (provides cascade of events), Gene Ontology, Ingenuity (paid), Biocarta (functional relationships between proteins involved), String Database, Human Protein Interaction Database; or you may create a manually curated by experts lists<br><br>    Public databases - Biased (not all genes are included), disease genes tend to be investigated more often, genes that are more investigated will have more interactions), not always reliable (interactions are often predicted, not validated; if experimentally tested, it is unknown how reliable it is) |
| **What was the problem when** | Comparing public databases and manually curated<br>    SPL list - Synapses |

| | |
|---|---|
| **the SPL list (manually curated) was compared to public databases?** | 438/1043 of SPL genes have no KEGG pathway<br>388/1043 of SPL genes have no gene ontology<br>655 are not synaptic according to the database |
| **On which aspects do tools for genes-set analysis differ?** | Tools for gene-set analysis<br>Differ in:<br>• Self-contained or competitive tests<br>• Different statistical algorithms test different alternative hypothesis<br>• Different sensitivity for LD, number of SNPs, number of genes, background heritability |
| **What is the difference between self-contained and competitive tests?** | Self-contained - H0: Genes in gene-set are not associated with traits; problematic because polygenic traits will always be associated with a large part of the genome<br>Competitive tests - H0: Genes in gene-set are not more strongly associated than genes **NOT** in the gene-set |
| **What happens with type one erros when you use a competitive test? And when you use a self-contained test?** | The more polygenic your trait is, the more likely you are to get a false positive when using self-contained tests (higher background heritability -> effect of the entire genome on the trait)<br><br>This doesn't occur with competitive tests<br>Never use self-contained tests |
| **What is the difference between minimal P-value and combined p-value?** | Different statistical algorithms test different alternative hypothesis<br>Minimal P-value - At least one SNP needs to be significant<br>Combined P-value - It doesn't matter if SNPs are signicant |
| **What happens with type 1 error when you** | Effect of gene size for signifance value |

| | |
|---|---|
| **increase gene size?** | **Gene size**<br><br><br><br>MAGENTA should not be used |
| **What happens with type 1 error when you increase LD?** | LD for significance value<br><br>**f   Linkage disequilibrium between genes**<br><br><br><br>FORGE/ALLIGATOR should not be used |
| **What happens with type 1 error when you increase number of genes?** | Number of genes for significance value<br><br> |

# 6b. MAGMA

| | |
|---|---|
| **What is MAGMA?** | MAGMA - Software tool for gene and gene-set analysis<br>    Input - Genotype and phenotype data; summary statistics for GWAS + gene definition and gene sets |
| **What are the three main steps of using MAGMA?** | Three main steps<br>    Annotation - Map SNPs onto genes<br>    Gene analysis - Compute association of genes with phenotype<br>    Gene-set analysis - Gene association in gene sets<br>        Extension modeling options: Advanced analysis (conditional and joint analysis) |
| **What is the annotation stage based on?** | Annotation - Based on physical location of SNPs/genes in the DNA<br>    Window around gene - Variable, mainly used for upstream/transcription start site of gene<br>    Genes can be non-protein-coding as well |
| **What are the three models for SNP association? Which one is the most recommended?** | Gene analysis<br>    Joint association of all SNPs in a gene with the phenotype<br>    Different analysis models have greater sensitivity to different genetic architectures<br>        MAGMA has different models with different sensitivies<br>    • Principal component linear regression - Requires raw genotype data<br>    • SNP-wise mean - Mean SNP association<br>    • SNP-wise Top - Strongest SNP association<br>    • SNP-wise Multi - Combines SNP-wise mean and top |
| **Why does the gene-set analysis use a one-sided statistical test?** | Gene-set analysis - Essentially a t-test<br>    Unit of analysis: Genes<br>    One-sided test - High negative association scores are not interesting<br><br>    Self-contained analysis (single sample t-test) - H0 - $\mu s = 0$<br>    Competitive analysis (dual sample t-test) - H0 - $\mu s = \mu 0$ |
| **What are some stastical challenges you might face when using MAGMA?** | Statistical challenges<br>    Outlier effects - Small subset of strongly associated genes driving the association<br>    Independent observations - Linkage disequilibrium<br>        Model covarience<br>    Confounding - Apparent effect may be induced by an overlap with truly associated set |

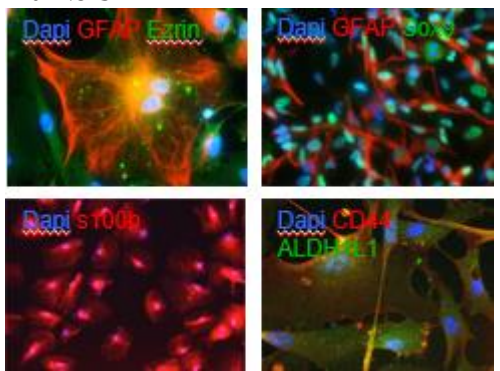| | |
|---|---|
| | Confounding factors can only be ruled out by experimental design |
| **Which statistical test is used in the competitive analysis of MAGMA?** | Competitive analysis<br>    Linear regression model - Two sample t-test (much more flexible)<br>        Example: Brain-specific gene expression<br>    Lower p = higher Z |
| **How can you detect an effect when a gene has a correlated expression across tissues?** | Correlated expression across tissues<br>    Solved by multiple regression - If the effect is true, the apparent effect of correlated variables will disappear<br>    &bull; Requires true effect variable to be present<br>    &bull; Very large number of variables (reduces power)<br>    &bull; This does not rule out confounding effects |
| **Why might you want to divide your gene set into different subsets?** | Interaction analysis<br>    Multiple linear regression - Analyse the overlapping of pathways<br><br>    Example: Division of miRNA-145 gene set into four subsets: Three non-significant, one very significant -> This effect would be undetectable by regular gene-set analysis |

# 7a. Modeling genetically complex disorders using induced pluripotent stem cells

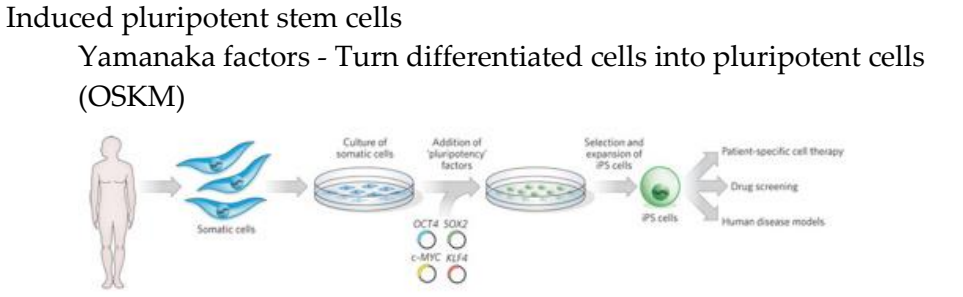| | |
|---|---|
| **What is the definition of a stem cell?** | Stem cells - Can divide or differentiate into different cell types  |
| **What do trophoblasts and blastocyst form?** | Trophoblast - Forms placenta <br> Blastocyst - Forms the person  |
| **What are the types of stem cell potency?** | Types of potency <br>     Totipotent - Can make an entire human being (extraembryonic/placental cells + all tissues) <br>       Before 16-morula stage <br>     Pluripotent - Can make all tissues (but cannot form an individual/form placental tissue) <br>     Multipotent - Can make more than one cell type, but not all (haematopoietic stem cells) |

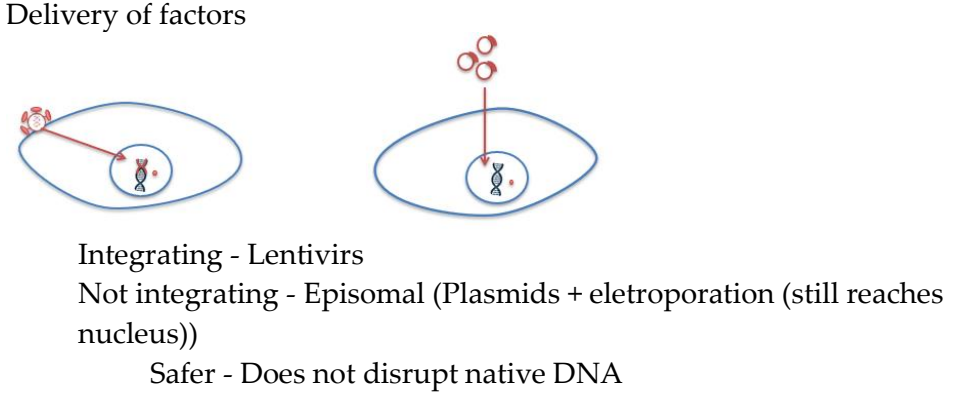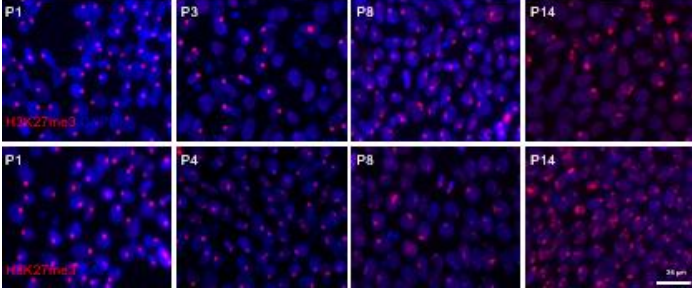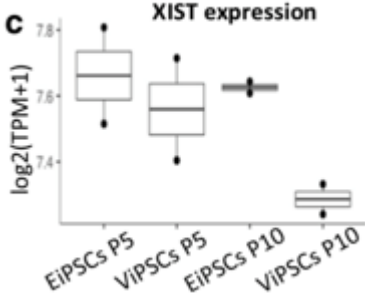| | |
|---|---|
| **How would you know if your stem cells have differentiated into the tissue you wanted?** | Model tissue in vitro<br>　Embryonic stem cells from the inner cell mass (pluripotent)<br>　Differentiation can be mimicked in vitro - Inserting the same factors in vivo<br><br>　How to know you have made the right tissue - Antibodies for cell markers<br> |
| **What are the four Yamanaka factors? How were they discovered?** | Induced pluripotent stem cells<br>　Yamanaka factors - Turn differentiated cells into pluripotent cells (OSKM)<br><br>　Process of discovery: 100 factors genes expressed in PSC -> narrowed to 24 (transcription factors since they need to modify expression of other genes); leave factors off one at a time |
| **How do Yamanaka factors work?** | Yamanaka factors:<br>　Somatic genes are silenced, pluripotent genes will be switched on<br>　Chromatin remodeling<br>　Epigenetic remodeling |
| **What are two options to deliver stem cell gene factors to a cell?** | Delivery of factors<br><br>　Integrating - Lentivirs<br>　Not integrating - Episomal (Plasmids + eletroporation (still reaches nucleus))<br>　　Safer - Does not disrupt native DNA |

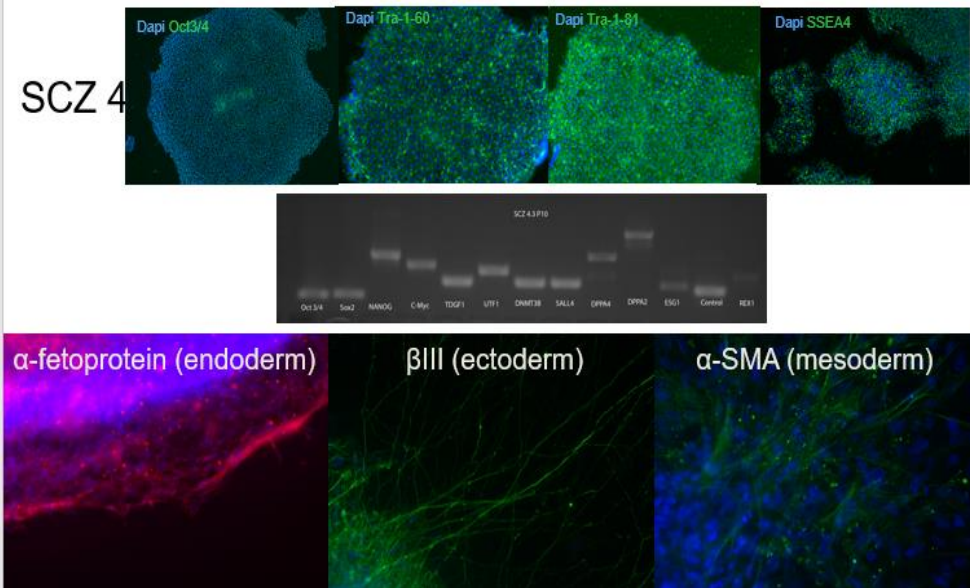| | |
|---|---|
| **What causes X-inactivation?** | Females are mosaics - Different inactivation of X chromosome<br>    XIST - RNA transcript activates PCR2 (adds methylation to DNA)<br>    H3K27 - Repressive chromatin marker of inactivation |
| **Which method of stem cells delivery factors would be safer for ex-vivo therapy for women? Why?** | <br>Inducing pluripotency removes X chromosome inactivation using Lentivirus (first row)<br>X chromosome inactivation is maintained when using plasmid (second row) |
| **What happens to XIST expression during lentivirus treatment?** | XIST expression remains normal in episomal reprogramming, not in lentivirus reprogramming<br><br>X-linked disorders could be treated with X chromosome re-activation of lentivirus |
| **Describe a pluripotency test. What is the X and Y axis? What you expect if a cell was converted successfully into PSC?** | Pluripotency test<br>    Ratio of gene expression of PSC (Y axis) compared to differentiated cells (X axis)<br><br>There is no difference between lentivirus and episomal |

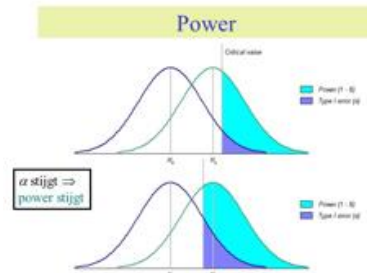| How would you prove that a cell line is pluripotent? | Prove that a cell line is pluripotent:<br><br>Immunostaining - Check for expression of PSC factor<br>Immunoblot - Expression of these protein<br>Three markers for the the three germ layers |
|---|---|
| Why should we use iPSC? | Benefits of iPSC<br>    Ethical - No embryos needed<br>    IPSCs are identical to donor - Safer transplantation, good for genetic disorders that are very complex (e.g. schizophrenia with 108 associated SNPs) |
| **Define the main differences between CNVs and SNPs associated with schizophrenia.** | Genetics of SCZ<br>    Heritability: 81%<br>    Lifetime risk: 1%<br><br>    CNVs (copy number variants): Rare, large effect, few patients, low penetrance<br>    SNPs: Common, very small effects (combined large effects), majority of patients |
| **How has been SCZ been studied in the past?** | Studying SCZ<br>    Post-Morten studies<br>    Animal studies: Pharmacology, transgenic mice (pick CNV or CRISPR 108 -> incredibly laborious)<br><br>    Picking CNV - DISC1, miR-137, NRXN1 -> Create synaptic deficit (they are all synpase proteins); creates a biased sample |

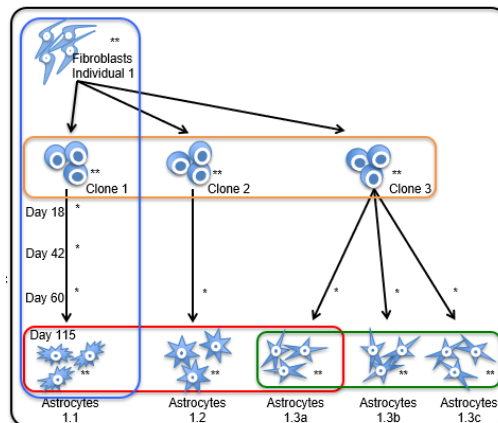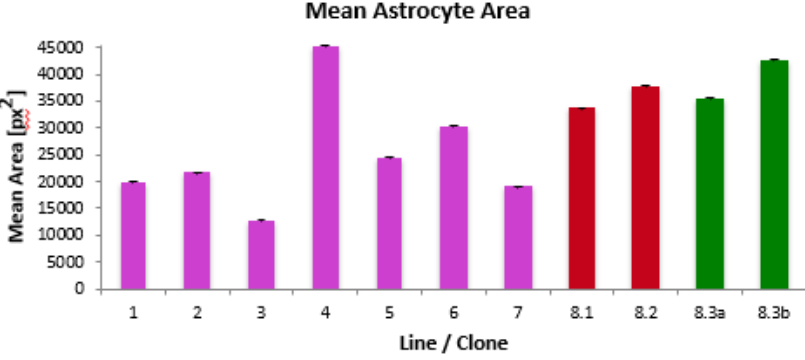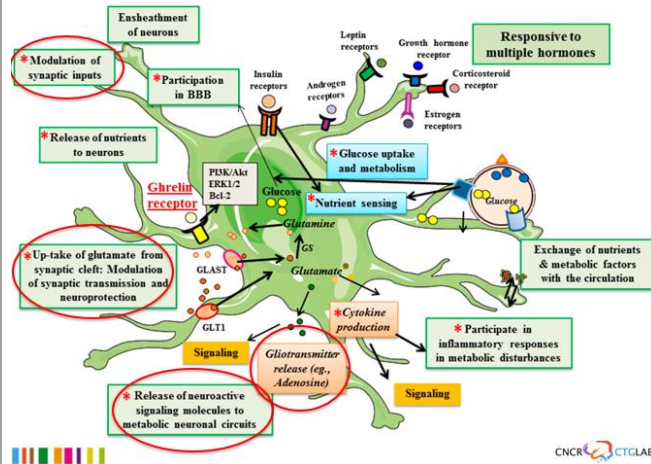| | |
|---|---|
| **How could you improve the power of your iPSC study?** | IPSC for SCZ<br>    Heterogeneous population -> decreases power<br>    Pick extreme SNPs from polygenic risk score<br>        Pull graphs apart - Controls Low-risk score and patients with high-risk score<br><br><br><br>    Problems: Expensive, laborious, various protocols (variability between studies)<br>        Accumulation of mutation<br>        Defects in reprogramming<br>        Inaccuracy of protocol |
| **Describe how you could validate your protocol to convert differentiated cells into iPSC.** | Determine variability and source<br><br><br><br>    Three individuals - 3 cell lines from each; analyse variability of cells before and after differentiation from the same individual<br>    Whole-proteome from single cells - Correlation of each cell to one another<br>        Differentiated cells are more similar to one another than to other people - Means that protocol is reliable; undifferentiated cells are not<br>            Older colonies become more alike - Hard to interpret<br>            Technical replication - Same sample multiple times |

| | |
|---|---|
| | Cellomics<br>    Astrocyte area - Variation is larger between individuals (pink) compared to clones from the same individual (red) or technical replicates (green)<br><br><br><br>**Mean Astrocyte Area** |
| **How many clones should we use per individual to optimize resources?** | Conclusions<br>    1 clone per individual and multiple people<br>    Sequence entire genome from each cell - Check for unexpected mutations (PRS, CNVs) |
| **Why should we not use lentivirus to study rare diseases?** | Episomal reprogramming is better for rare disease - Lentivirus requires integration, may disrupt DNA in unexpected ways |
| **What are some neuronal dysfunctions of SCZ patients?** | Neurons in SCZ<br>    Dysregulation of excitation-inhibition network<br>    Synaptic pruning -> brain volume decreases<br><br>    SMAD inhibition - Block pathways of differentiated cells<br>        Forms rosettes -> neural tube-like, makes both neurons and glial cells<br>    HSHH and Valproic acid - GABAergic cells<br><br>    Neuritis length is reduced in SCZ patients<br>    Dendrite length stays the same<br>    Axonal length is reduced in SCZ patients |

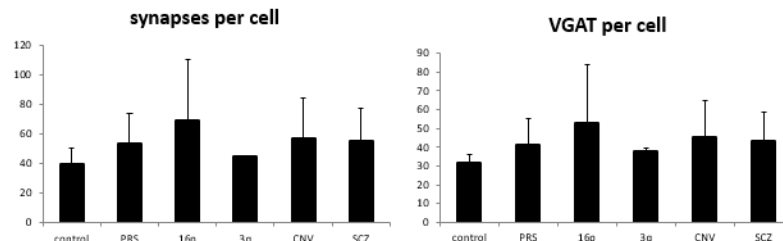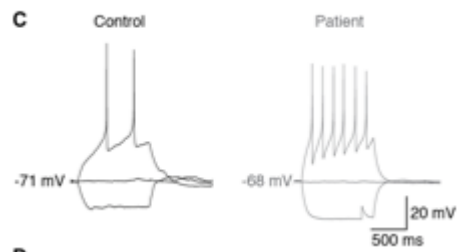| | |
|---|---|
| **What happens with astrocytes in SCZ?** | Astrocytes in SCZ<br><br><br><br>31 astrocyte specific gene sets in SCZ<br>Coexpressed genes (neurons and astrocytes) - CACNA1C<br>Less density of astrocytes and GFAP in SCZ patients |
| **What happens when SCZ astrocytes were co-cultured with neurons?** | Results SCZ<br>  Nestin/CD44 - Expressed in stressed cells<br>    No difference in astrocytes alone<br>When co-culturing SCZ astrocytes with neurons - VGAT (gabaergic terminals) is lower for SCZ patients; ratio of GABAergic and glutamatergic neurons is altered<br><br><br><br>  Sandwich culture - Cells not in physical contact, but astrocytes can secrete factors |
| **What happens with oligodendrocytes in SCZ?** | Oligodendrocytes in SCZ<br>  White matter abnormalities in patients<br>  Less MBP in chimeras (human iPSC in mice) |
| **Why should we use organoids to study SCZ?** | Organoids - three dimensional -> allows myelination to occur<br>  Reelin - Produced in the outer layer of the brain<br>  Ctip2 - Layer 5/6<br>  MBP - Myelinating Oligodendrocytes |

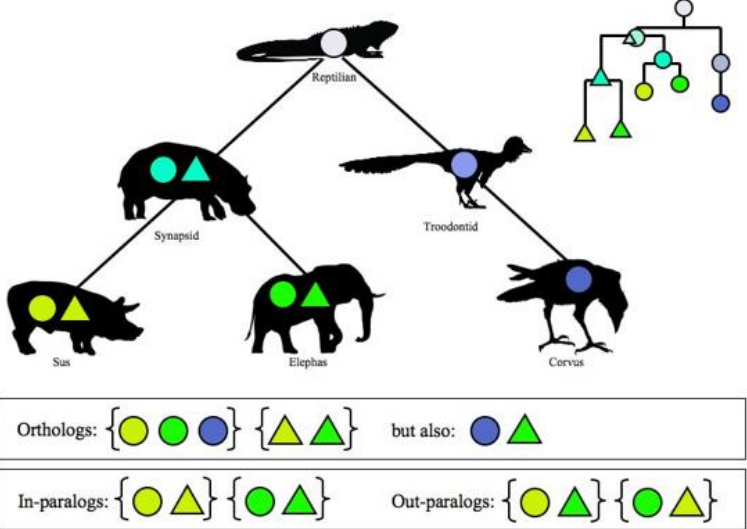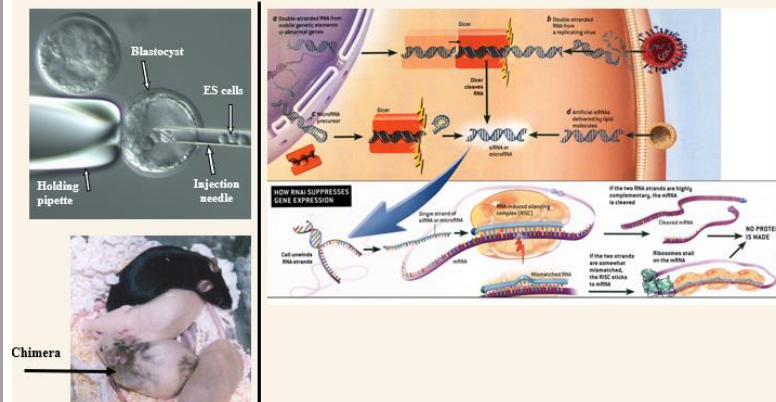| | |
|---|---|
| **What happens to the half-width of action potentials in SCZ patients? What does that mean?** | C    Control                    Patient<br><br>-71 mV ─                    -68 mV ─<br><br>                                                    ⎸20 mV<br>                                                    500 ms<br>D<br><br>Half-width of action potential - How long does a neuron take to go back to resting state<br>  Measure from depolarization and hyperpolarization peaks<br>  Half-width is longer for SCZ patients - It takes longer to recover |

# 7b. Gene function analysis in Neuroscience for Mendelian Disorders

| | |
|---|---|
| **How could we use yeast to study a human disease, like Wilson's?** | Wilson's disease - ATP7B<br>    Monogenic disease with a known gene - The observed allelic variant of the patient will or will not explain the disease<br>    ATP7B structure - If the mutation is on the ligand site, you may argue that the mutation is causal<br><br>    Wilson's disease in yeast - ccc2 gene from yeast is similar to ATPB7<br>        Copper transport -> important for iron metabolism<br><br>        Reinserting ATP7B in yeast restores function in knockout yeast<br>            Tests for different patient mutation<br>            Cheap, ethical and quick (three days) |
| **What are homologs, orthologs and paralogs?** | Homologs - Genes shared from a common ancestor; divided in ortholog and paralogs<br><br><br><br>    Orthologs - Similar structure and different function<br>    Paralogs- Genes related by duplication in a genome - have different function |
| **What can you do with a knockout?** | Knockout functions<br>    Investigate gene function<br>    Test causality of candidate disease mutation<br>    Develop and test therapies in disease model |

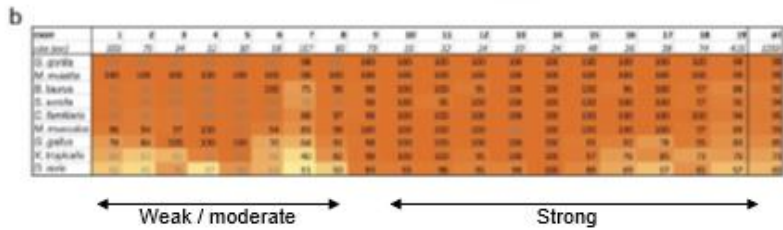| | |
|---|---|
| **What are two common methods to inactive genes?** | Methods to inactive genes  Gene knockout<br>RNA interference - Translation will be blocked or mRNA will be degraded very quickly<br>　　Very efficient in zebrafish - Double stranded modified RNA to be more stable (morpholinos) |
| **What is the phenotype of someone with AUTS2 deletion?** | AUTS2 in a syndromes form of autism<br>　　AUTS2 deletion - Autistic behavior, IQ decline, short stature, microencephaly, feeding difficulties,generalized hypotonia (muscle weakness)<br>*Heterozygous loss of gene is sufficient to cause the phenotype* |
| **What correlates with the severity of phenotype of AUTS2-deletion patients?** | All mutations were intronic - Variant of unknown significance<br>　　AUTS2 deletions are observed in patients with these symptoms<br>　　Heterogeneity of phenotypes - Sum of presence of scores = AUTS2 score<br><br>　　　　Phenotypes are more severe in patients with more deletions in the C-terminus - Unexpected because N-terminus truncation should be more severe |

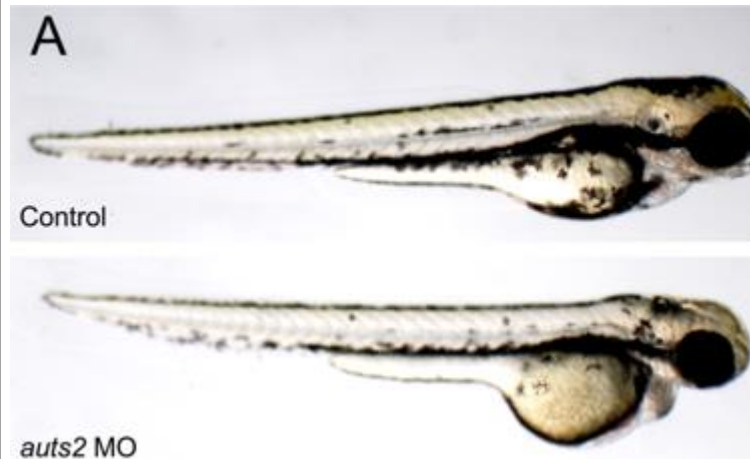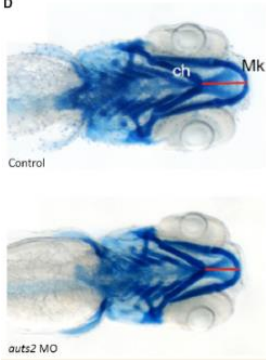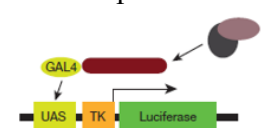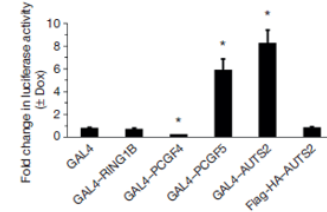| | |
|---|---|
| **What is the explanation for the fact the deletions in the C-terminus are more severe?** | Evolutionary conservation of AUTS2<br><br>    Amino acid identity - C-terminus is more conserved evolutionarily than N-terminus<br><br><br><br>    Zebrafish may have a functional ortholog to the human protein |
| **What method could you use to identify an alternative transcriptional splice site?** | Alternative transcriptional splice site<br>    Present in the middle of exon 9 - There is two gene products inside the same gene<br>    RACE - Rapid Identification of cDNA Ends |
| **How could we use zebrafish to study a human deletion like AUTS2?** | Zebrafish<br><br><br><br>    Morpholinos - Knockdown of AUTS2 causes microencephaly (establishes causality between gene and phenotype)<br>    Phenotype is rescued by the insertion of human AUTS2 cDNA into zebrafish<br>        cDNA does not have introns - It is easier to produce in silico<br>    Microencephaly is also rescued by the C-terminal short isoform (cDNA) |
| **How can facial dysmorphism be quantified in zebrafish?** | Facial dysmorphism |

| | |
|---|---|
| | Blue - Stain for cartilage<br>Larger distance between Ch and Mk (ceratohyal and Meckel's) |
| **What happens to the number of neurons in AUTS2 zebrafish?** | Reduced proliferation of neuronal progenitors in auts2 morphants<br>    Reduction of differentiated neurons<br>    Phosphohistone-h3 antibody - Epigenetic markers (proliferating neurons)<br>        There are less proliferating precursor cells |
| **What is the main advantage of studying mendelian disorder, as opposed to complex traits?** | Mendelian disorders - Sample size can be much smaller<br>    Causal role of mutation - Large effect size |
| **What is the mechanistic function of AUTS2? Why is it unusual?** | AUTS2 gene product function - Still unknown<br>    Possible transcription factors<br>    Enhancer regions in the middle of AUTS2 gene<br>    AUTS2 is part of Polycomb complex - Positively regulates transcription of targeted genes (while most other components of this complex reduce transcription)<br>AUTS2 binds to promoter regions<br>Mouse models confirm the finding |