# Lecture 8: Genetic variants and their interpretation
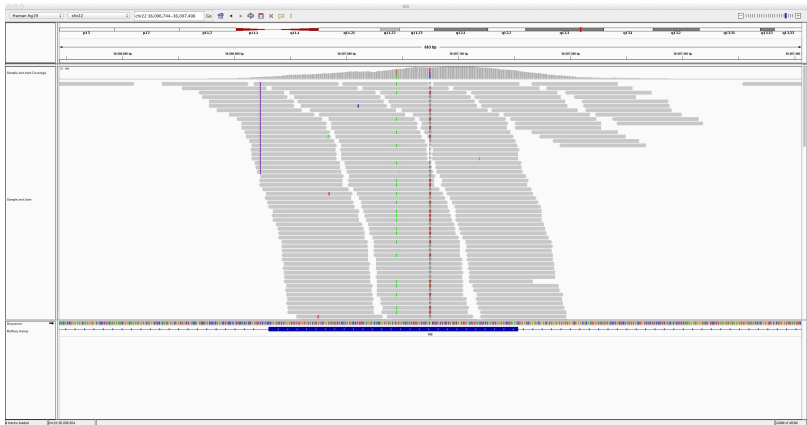
Paolo Provero

# NGS and genetic variants

NGS techniques allow us to detect the *genetic variants* that characterize each individual

- *Targeted sequencing*: detect variants in selected regions
  - Whole Exome Sequencing
  - Sequencing of candidate genes or exons
  - The smaller the target the greater the sequencing *depth* with the same number of reads
- *Whole genome sequencing*

# Variant detection

# VCF format

To save space, the variants are expressed as *differences* from the reference sequence

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS      ID        REF ALT    QUAL FILTER INFO                              FORMAT      NA00001        NA00002        NA00003
20     14370    rs6054257 G   A      29   PASS   NS=3;DP=14;AF=0.5;DB;H2           GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
20     17330    .         T   A      3    q10    NS=3;DP=11;AF=0.017              GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3   0/0:41:3
20     1110696  rs6040355 A   G,T    67   PASS   NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2   2/2:35:4
20     1230237  .         T   .      47   PASS   NS=3;DP=13;AA=T                  GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20     1234567  microsat1 GTCT G,GTACT 50  PASS   NS=3;DP=9;AA=G                   GT:GQ:DP    0/1:35:4       0/2:17:2       1/1:40:3
```

# Variant interpretation

Problem: determine which variants are causal of a phenotype/disease

- *Mendelian* phenotypes
    - *rare* variants
    - high *penetrance*
    - in coding regions
- *Complex* phenotypes
    - *common* variants
    - non-coding regions

# Effect on the phenotype

The likelihood that a variant will affect the phenotype depends on where the variant is located (coding exon, UTR, intron, intergenic)

- Coding variants can be classified based on their effect on the protein sequence
  - synonymous: no change in protein sequence (possible due to the degenerate nature of the genetic code)
  - missense: changes one aminoacid into another
    - conservative: new aa is chemically similar to old one
    - non-conservative: new aa is chemically different – likely change in structure
  - nonsense: introduces a stop codon
- Non-coding variants are assumed to affect phenotype by changing gene regulation
  - e.g. by creating/destroying a transcription factor binding site
  - much more difficult to classify/interpret

# Coding variants: the XLMR example



*Nature Genetics 2009*

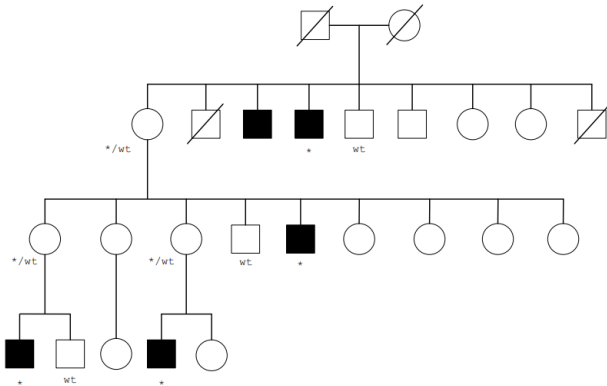# Coding variants: the XLMR example

- 208 families with MR and pattern of transmission compatible with X linkage
- Exome sequencing of ~700 genes in X chromosome
- Most differences in coding sequence between individuals are recurring and found in dbSNP
- Consider *truncating* variants only (found in 30 genes)

# Segregation

A variant can be *causal* of the disease only if it *segregates* - presence of the variant in all affected members and only in affected members
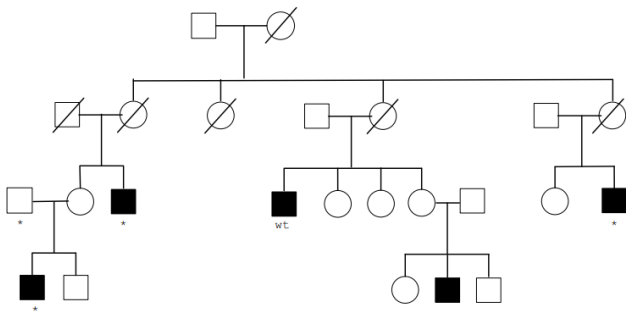


Family 62
UPF3B(NM_080632.1) c.1288C>T p. R430X

# Non-segregation

**Family 74**
ITIH5L (NM_198510.1) IVS7+1ins T

# Message

- Exome sequencing can help finding causal variants of Mendelian diseases
- However even truncating variants can be compatible with normal phenotype
  - "loss of function of 1% of the genes in the X chromosome is compatible with apparently normal existence"
- Another recent study found an average of ~40 homozygous LOF mutations in normal individuals
- Therefore when sequencing the exome of a *proband* with a genetic disease we expect to find many candidate mutations
  - Need for *variant prioritization*
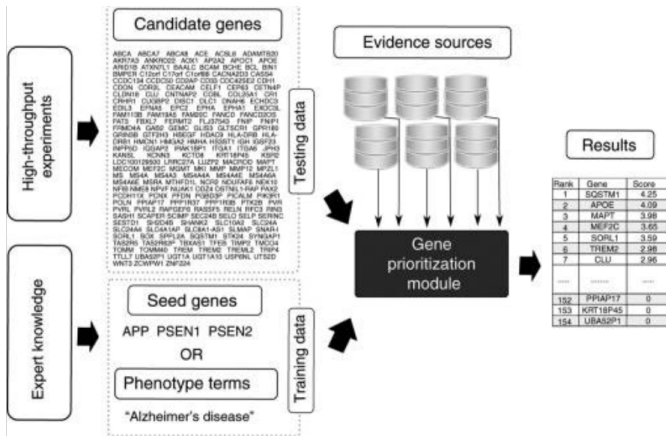
# Variant prioritization

Problem:

- Among all the variants found in a proband with a genetic disease, find the ones most likely to be causative

The following criteria increase the probability that a variant is causative

- very low frequency in the population
- predicted strong effect on the protein
- gene involved in biological process relevant to the disease
- gene interacts with genes involved in the same/similar diseases
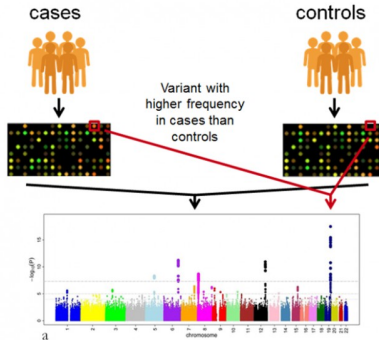
# Variant prioritization



*From Zolotareva & Kleine, J Integr Bioinform 16:20180069 (2019).*

# Complex phenotypes

- Complex trait: determined by both genetic and non-genetic (e.g. environmental, behavioral, . . . ) factors
  - Diabetes
  - Rheumatoid arthritis
  - Crohn disease
  - Height
  - Blood pressure
  - Lymphocyte count
- Genetic determinants are usually
  - Many variants of small effect
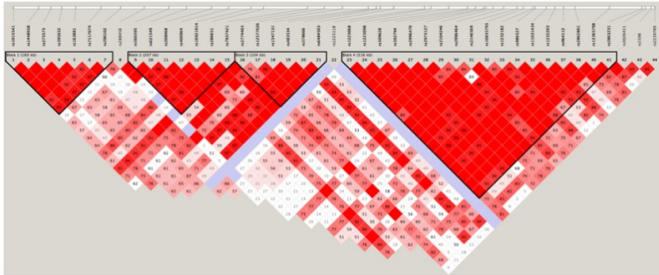  - Non-coding (hence presumably regulatory)

# GWAS

Genome-wide associations studies look for variants with significantly different frequency in cases and controls

# Predicting vs understanding complex traits

- GWAS hits can be used to build *polygenic scores* predicting disease risk
- To understand disease (hence possibly find cures) one needs to understand the mechanism leading from variant to disease
- GWAS hits are mostly non-coding $\rightarrow$ regulatory effect
  - What is the *target* of a regulatory variant (gene whose expression is altered by the variant)?
  - What is the effect of the variant?
    - Regulatory code much less understood than genetic code
  - Which one is the causal variant?
    - *Linkage disequilibrium* implies that many neighboring variants are inherited together, thus showing the same correlation with the disease although presumably only one or a few are causative
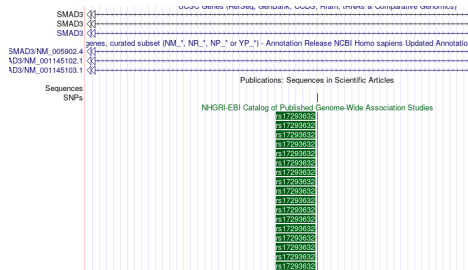
# Linkage disequilibrium



Strong coorelation between nearby SNPs indicate that they are
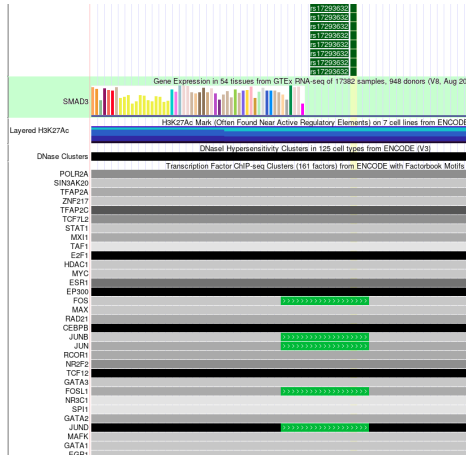inherited together (no recombination events)

# An example: rs17293632

- Associated by GWAS to many complex traits, including Crohn's disease
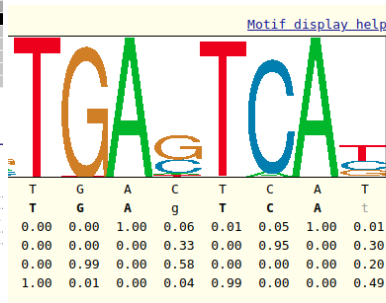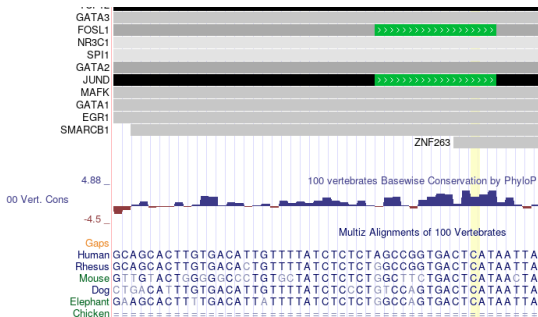- Located in an intron of SMAD3

# An example: rs17293632



- Colocalizes with several transcription factor binding peaks from ChIP-seq
- In particular, is within the core binding motif of AP-1

# An example: rs17293632

The variant disrupts a
highly conserved AP-1
motif

- The variant rs17293632 disrupts a highly conserved (hence, probably functional) AP-1 binding site
- Thus altering gene expression, presumably of SMAD3
- Thus altering the phenotype, e.g. by conferring susceptibility to Crohn's disease

# Genetic and epigenetic fine mapping of causal autoimmune disease variants

Kyle Kai-How Farh[1,2]*, Alexander Marson[3]*, Jiang Zhu[1,4,5,6], Markus Kleinewietfeld[1,7]†, William J. Housley[7], Samantha Beik[1], Noam Shoresh[1], Holly Whitton[1], Russell J. H. Ryan[1,5], Alexander A. Shishkin[1,8], Meital Hatan[1], Marlene J. Carrasco-Alfonso[9], Dita Mayer[9], C. John Luckey[9], Nikolaos A. Patsopoulos[1,10,11], Philip L. De Jager[1,10,11], Vijay K. Kuchroo[12], Charles B. Epstein[1], Mark J. Daly[1,2], David A. Hafler[1,7]§ & Bradley E. Bernstein[1,4,5,6]§

# eQTLs

To understand GWAS hits we need a systematic analysis of the effect of variants on gene expression.

An eQTL (expression quantitative trait locus) is a variant significantly associated with the expression of a gene
- gene expression considered as a *quantitative trait*, like height or lymphocyte count

An eQTL study needs a large cohort (hundreds of individuals) for which we have

- expression data
- genotyping (i.e. *dosage* [0,1, or 2] of each genetic variant)

The analysis is performed by *linear regression*

# Regression

- *Regression* is the main statistical method for the study of the *dependence* among two (or more) variables $x$ and $y$
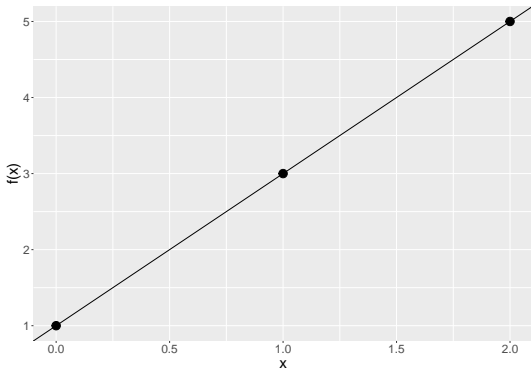- Assumes the existence of a *regression function* $f(x)$ which represents the true dependence of $y$ on $x$ as

$$y = f(x) + \epsilon$$

- In our case:
  - $x$: dosage of a variant ($x \in \{0, 1, 2\}$)
  - $y$: expression of a gene
- The *error term* $\epsilon$ represents random fluctuations, or the effects of variables that we do not consider
  - e.g. environmetal effects on the expression $y$ of the gene

# Linear regression

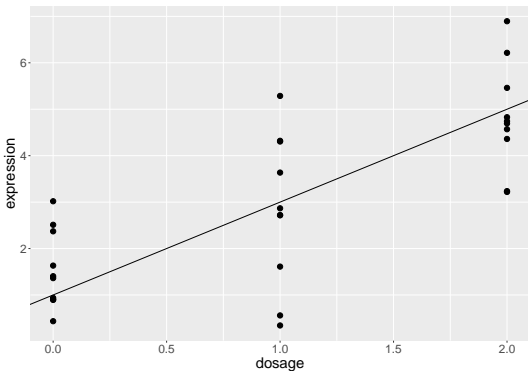In *linear regression* we assume $f(x)$ to be a straight line:

$$f(x) = \beta_0 + \beta_x x$$



In our case this means that we consider the effects of the paternal and maternal alleles on gene expression as independent and additive.

# Error term

If the regression function is linear we expect, considering the error term $\epsilon$, the actual measurements to look like this

## Estimating $\beta_0$ and $\beta_x$

Given the data:

- dosage $x_i$ ($i = 1 \ldots N$) and expression $y_i$ ($i = 1 \ldots N$) for a large number $N$ of individuals

our first goal is to estimate the regression function, i.e. the values of $\beta_0$ and $\beta_x$.

This is done by choosing the values that minimize the *mean square error*

$$MSE = \sum_{i=1}^{N}(y_i - f(x_i))^2 = \sum_{i=1}^{N}(y_i - \beta_0 - \beta_x x_i)^2$$
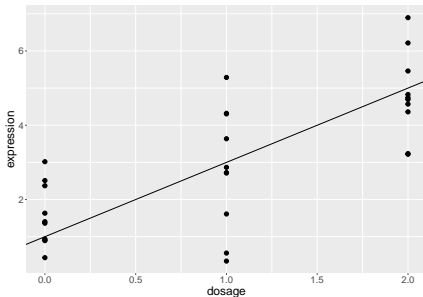
# Is this an eQTL?

Regression function:

$$f(x) = \beta_0 + \beta_x x$$

The variant is called an eQTL for the gene if $\beta_x \neq 0$ (if $\beta_x = 0$ expression does *not* actually depend on the variant dosage).

- Null hypothesis: $\beta_x = 0$
- Regression algorithms provide us not only with the value of $\beta_x$ (and $\beta_0$) but also with their *uncertainties*
- These can be used to *test* the null hypothesis and obtain a *P*-value
- Small *P*-values indicate that
    - the null hypothesis is probably false
    - that is, $\beta_x \neq 0$
    - that is, the variant dosage does indeed affect the expression of the gene
    - that is, the variant is an eQTL of the gene

# Example

For our (fake) data



- the estimated $\beta_x$ is 1.64
- with an uncertainty (standard error) of 0.277
- so that it is very unlikely that the true $\beta_x$ is 0
- indeed the $P$-value is $2.36 \cdot 10^{-6}$.

# The GTEx project

- eQTLs can be tissue-specific
- e.g. by altering the binding site of a brain-specific TF
- GTEx: eQTL analysis in ~50 human tissues
  - analyze all variants within 1 Mb of each gene

# Back to rs17293632

GTEx confirms that rs17293632 is an eQTL of SMAD3 in two tissues (but also of other neignboring genes)
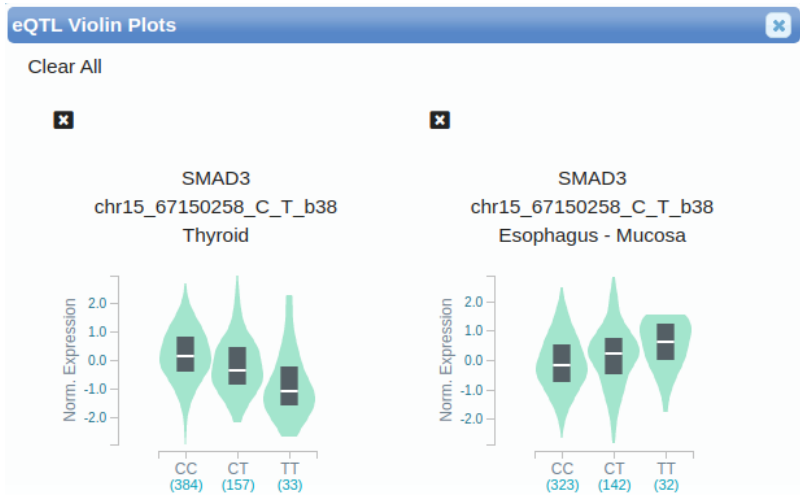


eQTLs of chr15_67150258_C_T_b38

Copy  CSV

| Gencode Id | Gene Symbol | Variant Id | SNP | | P-Value | NES ⓘ | Tissue |
|---|---|---|---|---|---|---|---|
| ENSG00000166949.15 | SMAD3 | chr15_67150258_C_T_b38 | rs17293632 | dbSNP ☑ | 5.8e-30 | -0.48 | Thyroid |
| ENSG00000166949.15 | SMAD3 | chr15_67150258_C_T_b38 | rs17293632 | dbSNP ☑ | 0.0000038 | 0.22 | Esophagus - Mucosa |
| ENSG00000103591.12 | AAGAB | chr15_67150258_C_T_b38 | rs17293632 | dbSNP ☑ | 0.000017 | -0.13 | Whole Blood |
| ENSG00000103591.12 | AAGAB | chr15_67150258_C_T_b38 | rs17293632 | dbSNP ☑ | 0.000064 | -0.10 | Esophagus - Mucosa |
| ENSG00000033800.13 | PIAS1 | chr15_67150258_C_T_b38 | rs17293632 | dbSNP ☑ | 0.00024 | 0.094 | Thyroid |

Showing 1 to 5 of 5 entries

# Tissue-dependent effects

The effect of the variant on SMAD3 expression depends on the tissue:

# GWAS and eQTLs

- eQTLs are useful in interpreting variants found by GWAS to be associated to complex traits/diseases
- Enrichment of eQTLs has been found among GWAS hits
- Recently a more global approach to the use of eQTLs has emerged: intermediate molecular phenotypes
    - In particular, *transcriptome-wide association studies* (TWAS)

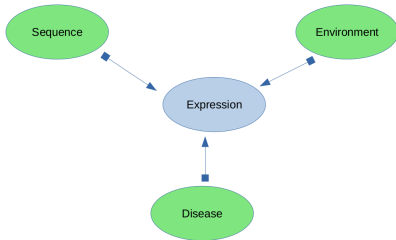# Gene expression as an intermediate molecular phenotypes

GWAS:



As most hits are non-coding, we actually believe



Gene expression is an *intermediate molecular phenotype*

# Components of gene expression



- Gene expression is not determined solely by genetics (sequence)
- However only the *genetic component of gene expression* (GREx) can mediate between variants and disease

# Transcriptome-wide association studies

**TECHNICAL REPORTS**

A gene-based association method for mapping traits
using reference transcriptome data

Eric R Gamazon[1,2,9], Heather E Wheeler[3,9], Kaanan P Shah[1,9], Sahar V Mozaffari[4], Keston Aquino-Michaels[1], Robert J Carroll[5], Anne E Eyler[6], Joshua C Denny[5], GTEx Consortium[7], Dan L Nicolae[1,4,8], Nancy J Cox[1,2,4] & Hae Kyung Im[1]

- Use eQTL data *from control individuals* (e.g. GTEx) to build a model predicting expression from genotype

- Use the models to compute the *GREx* of all genes for the GWAS diseased and control individuals

- Look for genes whose GREx correlates with the disease

- *Genetic variants affect the disease through the expression of these genes*

- The genes might be *therapeutic targets*

# Other molecular QTLs

## Genome-wide identification of DNA methylation QTLs in whole blood highlights pathways for cardiovascular disease

Tianxiao Huan ✉, Roby Joehanes, Ci Song, Fen Peng, Yichen Guo, Michael Mendelson, Chen Yao, Chunyu Liu, Jiantao Ma, Melissa Richard, Golareh Agha, Weihua Guan, Lynn M. Almli, Karen N. Conneely, Joshua Keefe, Shih-Jen Hwang, Andrew D. Johnson, Myriam Fornage, Liming Liang ✉ & Daniel Levy ✉

## Pooled ChIP-Seq Links Variation in Transcription Factor Binding to Complex Disease Risk

Ashley K. Tehranchi,[1] Marsha Myrthil,[2] Trevor Martin,[1] Brian L. Hie,[3] David Golan,[4,5] and Hunter B. Fraser[1,*]
[1]Department of Biology, Stanford University, Stanford, CA 94305, USA
[2]Department of Human Genetics, University of Chicago, Chicago, IL 60637, USA
[3]Department of Computer Science
[4]Department of Genetics
[5]Department of Statistics
Stanford University, Stanford, CA 94305, USA
*Correspondence: hbfraser@stanford.edu
http://dx.doi.org/10.1016/j.cell.2016.03.041

## A Genome-Wide Metabolic QTL Analysis in Europeans Implicates Two Loci Shaped by Recent Positive Selection
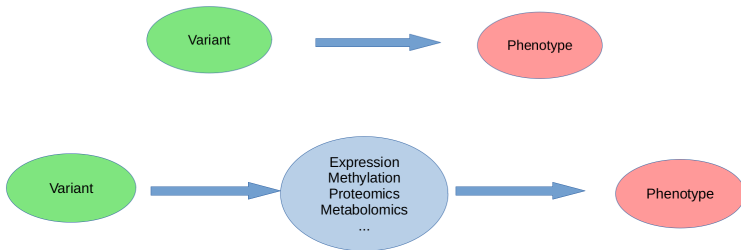
George Nicholson ✉, Mattias Rantalainen, Jia V. Li, Anthony D. Maher, Daniel Malmodin, Kourosh R. Ahmadi, Johan H. Faber, Amy Barrett, Josine L. Min, N. William Rayner, Henrik Toft, Maria Krestyaninova, Juris Viksna, [ ··· ], Chris C. Holmes ✉
[ view all ]

# Intermediate molecular phenotypes

# TWAS generalizations

Use other molecular phenotypes as intermediate phenotypes...

## PWAS: proteome-wide association study—linking genes and phenotypes by functional variation in proteins

Nadav Brandes[1]*, Nathan Linial[1] and Michal Linial[2]*

... or even macroscopic phenotypes

## Imaging-wide association study: Integrating imaging endophenotypes in GWAS

Zhiyuan Xu, Chong Wu, Wei Pan*, for the Alzheimer's Disease Neuroimaging Initiative[1]

Division of Biostatistics, School of Public Health, University of Minnesota, Minneapolis, MN 55455, USA

# Summary

Understanding the relationship between variants and phenotypes/diseases presents different challenges for Mendelian and complex phenotypes

- Mendelian diseases: find the causal coding variant through *variant prioritization*

- Complex diseases: understand the effect of many variants on gene regulation

In both cases the understanding of the causal relationship can in principle lead to new therapeutic strategies