# Identifying drug-targetable key drivers of disease

Expression data ——

—— Public data

Phenotypes ——

**UMCG**
Genetics Department

'To capture something small
you need something big'
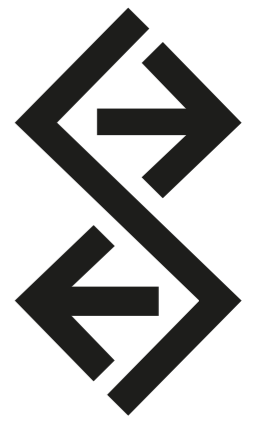
CERN

© Ruben van Leer

DNA ——————— AC
            CG
            GT

'To capture something small
you need something big'

DNA
Sequencers

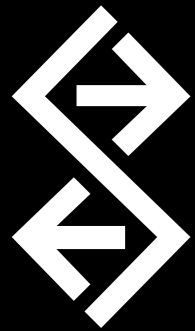'To capture something small you needed something big'

DNA Sequencer

# Minion
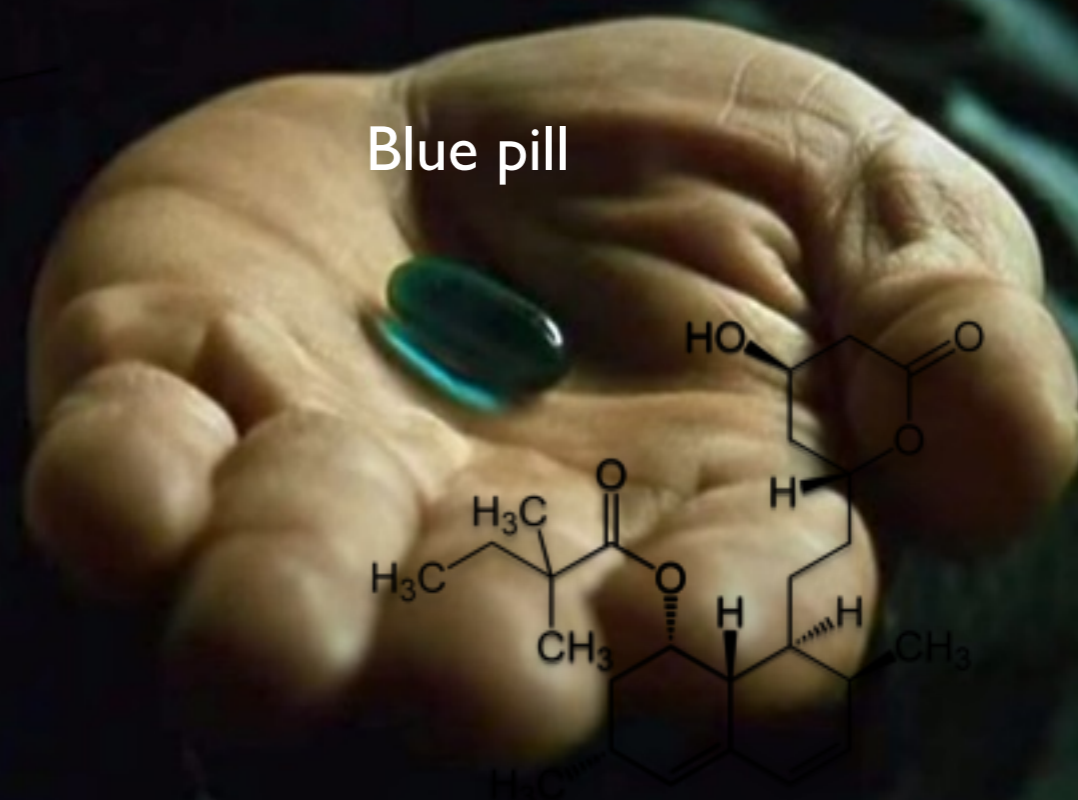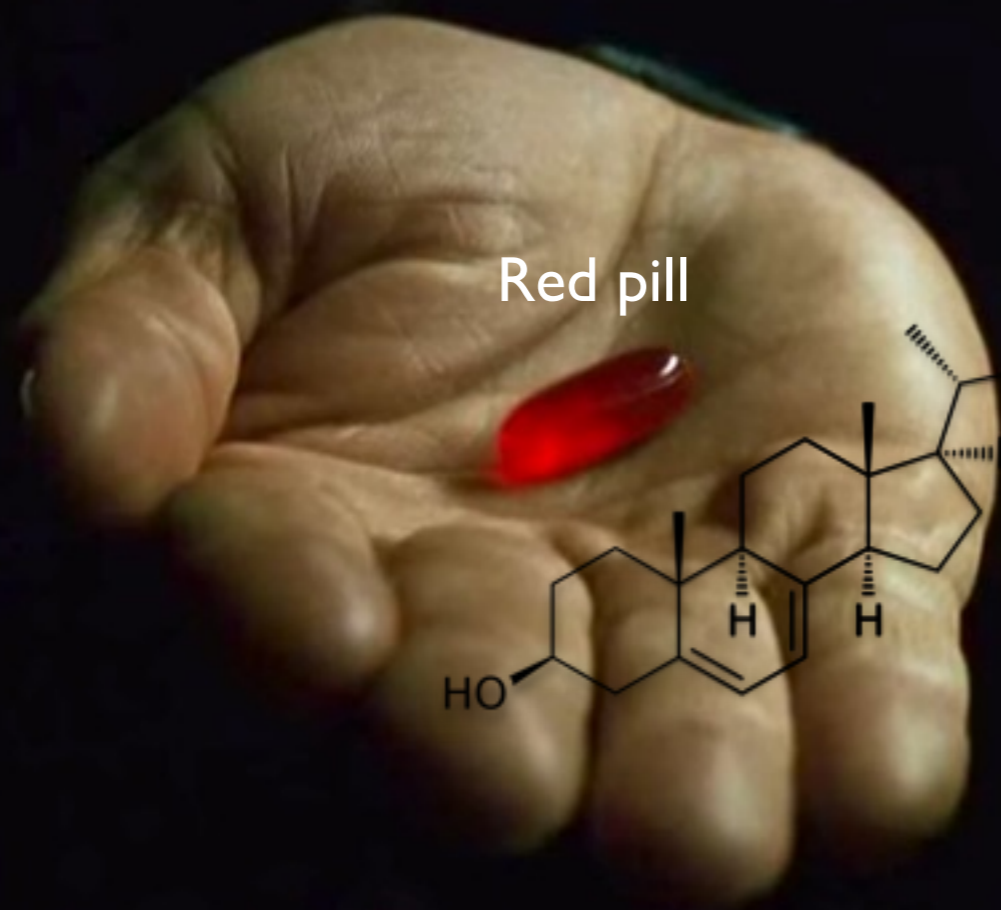
large amounts of data now available

Goal: better diagnose and treat patients

Red pill

Blue pill

6,054 disease associations

Age-related macular degeneration

2005

Teri Manolio *et al*: A catalog of published genome-wide association studies

## *Cis*-eQTL

SNP A/G

Gene X →

5′ —————————————— 3′

promoter region    exon 1    intron 1    exon 2

Gene X expression level

AA    AG    GG

## *Trans*-eQTL

Gene X    Protein X    Gene Y

Coding SNP A/G

5′ —————————————— 3′

promoter region    exon 1    intron 1    exon 2

**amino acid change**

5′ —————————————— 3′

promoter region    exon 1

Gene Y expression level

AA    AG    GG

Dubois *et al*, Nature Genetics 2010     Fu *et al*, PLoS Genetics 2012
Fehrmann *et al*, PLoS Genetics 2011     Westra *et al*, Nature Genetics 2013

Systemic lupus erythematosis risk factor:                    Chr. 7

Local expression effect:                    *IKZF1*                    Chr. 7

Type 1 interferon response:          *IFI6*   *IFI44L*   *IFIT1*   *MX1*          Downstream
                                                                                 *trans*-eQTL
(in Monocytes)                                                                   effects

Downstream effects identified for >200 genetic risk factors
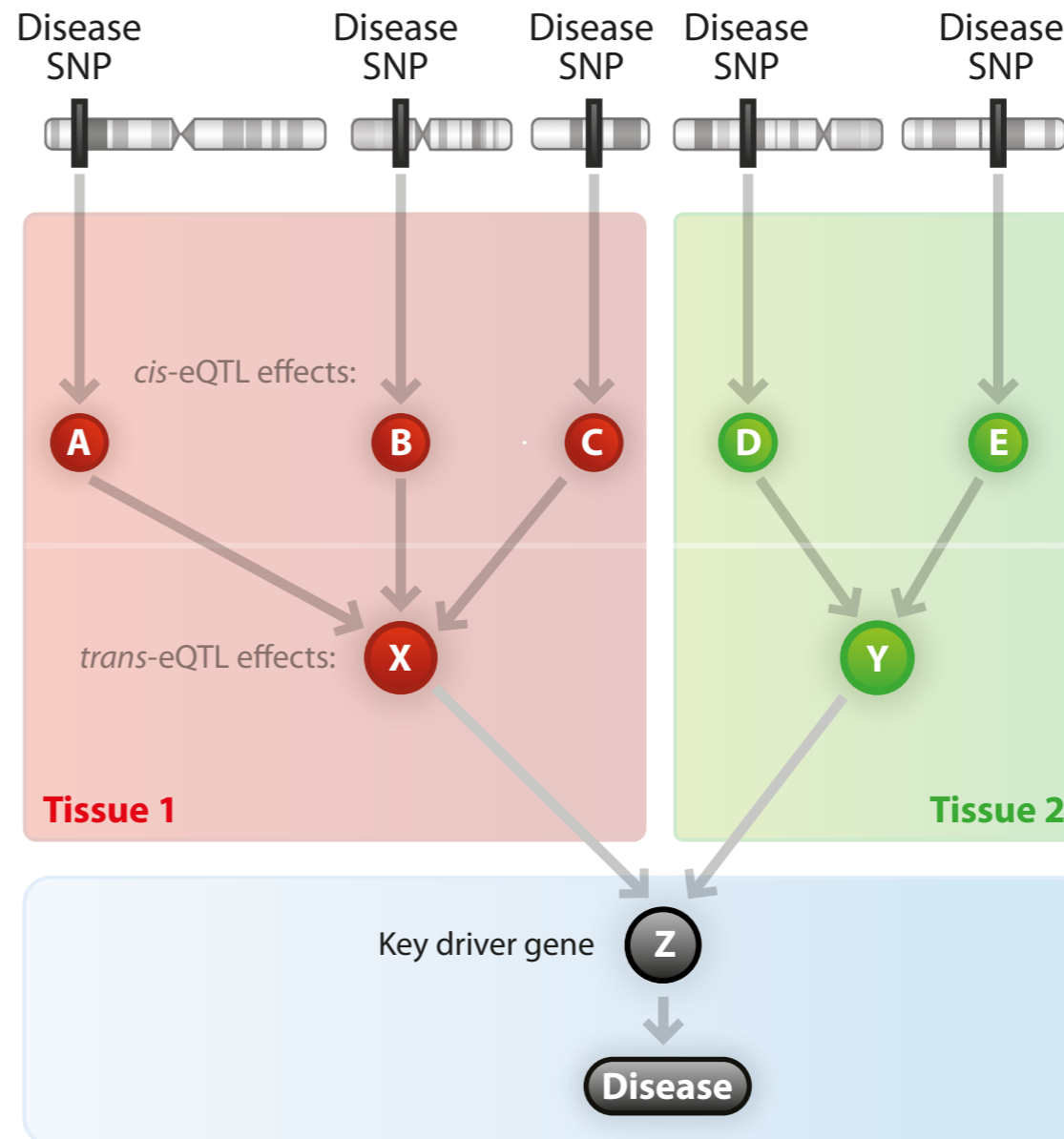
New meta-analysis ongoing in 25,000 blood samples

Genome-wide
association studies

*cis*-eQTL mapping

*trans*-eQTL mapping

Key driver gene
identification

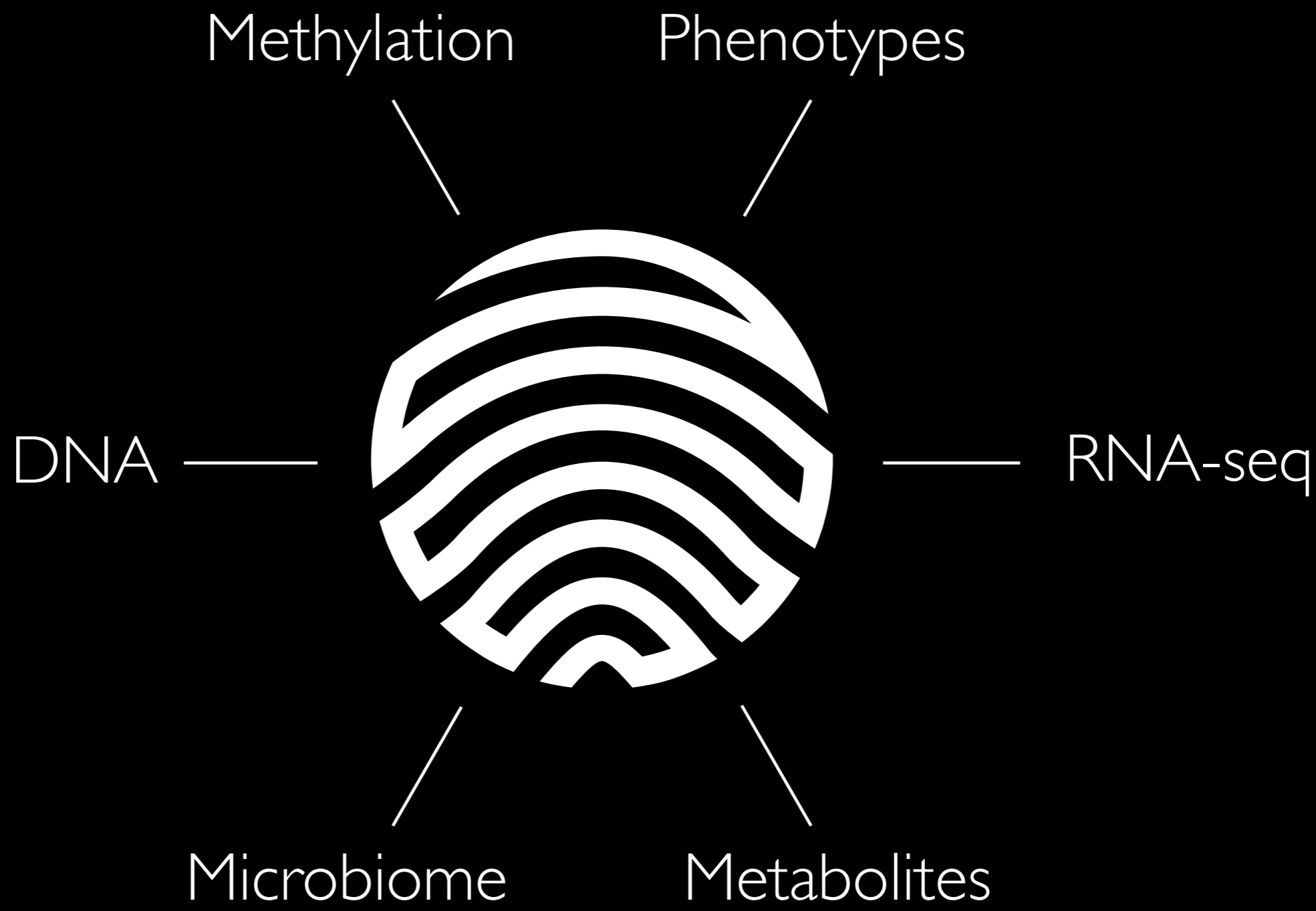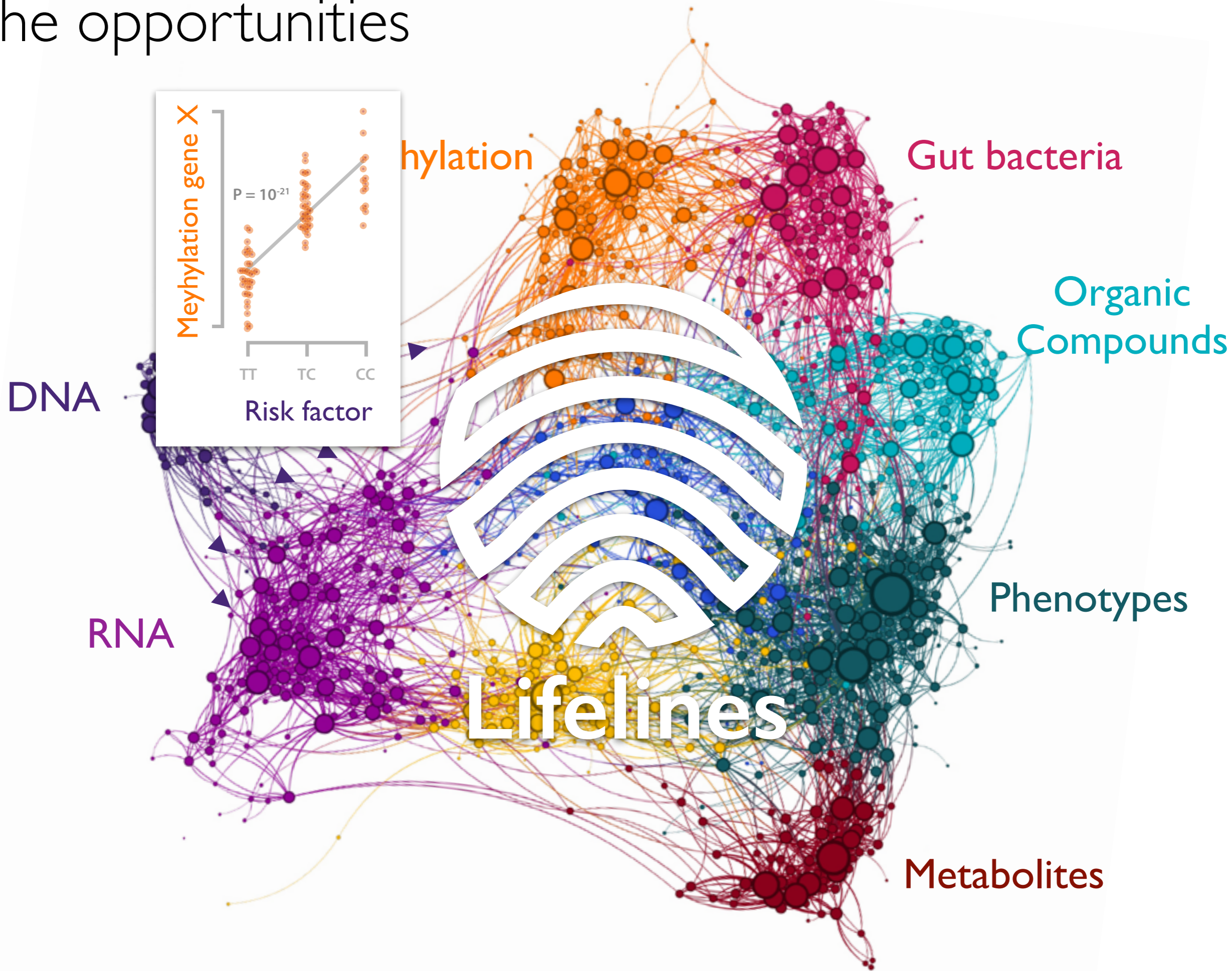# Impossible

**This is not going to be possible!**
- Massive sample-sizes required
- Many cell-types required
- Genotype and gene expression data required from the same samples

Methylation     Phenotypes

DNA ——

—— RNA-seq

Microbiome     Metabolites

**Lifelines Deep** (1500 samples)

# The opportunities



Meyhylation gene X

P = 10$^{-21}$

TT   TC   CC

Risk factor

...hylation

Gut bacteria

Organic Compounds

DNA

RNA

Phenotypes

Metabolites

Lifelines

- 34.4% of 405,709 tested CpG sites are *cis*-meQTL (FDR < 0.05)

- 31.2% of established GWAS risk factors give *trans*-meQTL effect (FDR < 0.05). 1,907 SNPs affecting 10,141 unique CpG sites in *trans*

- *Trans*-meQTL replicate in monocytes: 95% identical allelic direction

- *Trans*-SNPs affect expression of nearby TFs, subsequent methylation of downstream targets of these TF

**cis-eQTL**

Normalized expression

0.3
0.2
0.1
0.0
-0.1
-0.2
-0.3

$P = 2.6 \times 10^{-12}$

C/C      C/T      T/T

*NFKB1*
4: 103,422,486-103,538,459

**rs3774937**
4:103,434,253

**Risk factor associated to**
Ulcerative Colitis

**NOD2 eQTL in whole peripheral blood**

$P = 1.11 \times 10^{-294}$

NOD2

C/C    C/T    T/T

Leprosy risk SNP rs1981760

**NOD2 eQTL interaction analysis, STX3 interacts with rs1981760**

Interaction $P = 1.1 \times 10^{-69}$

C/C
C/T
T/T

NOD2

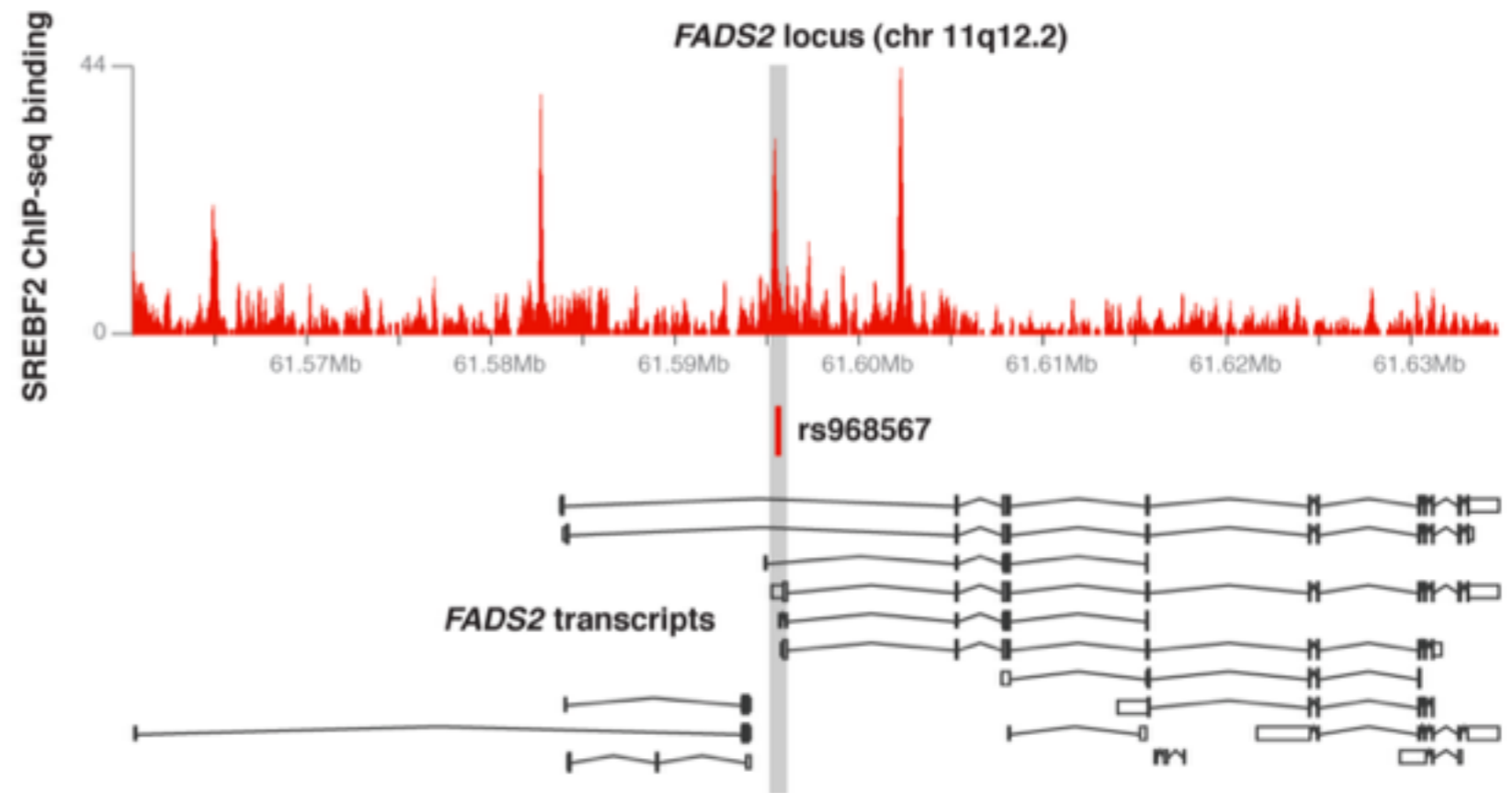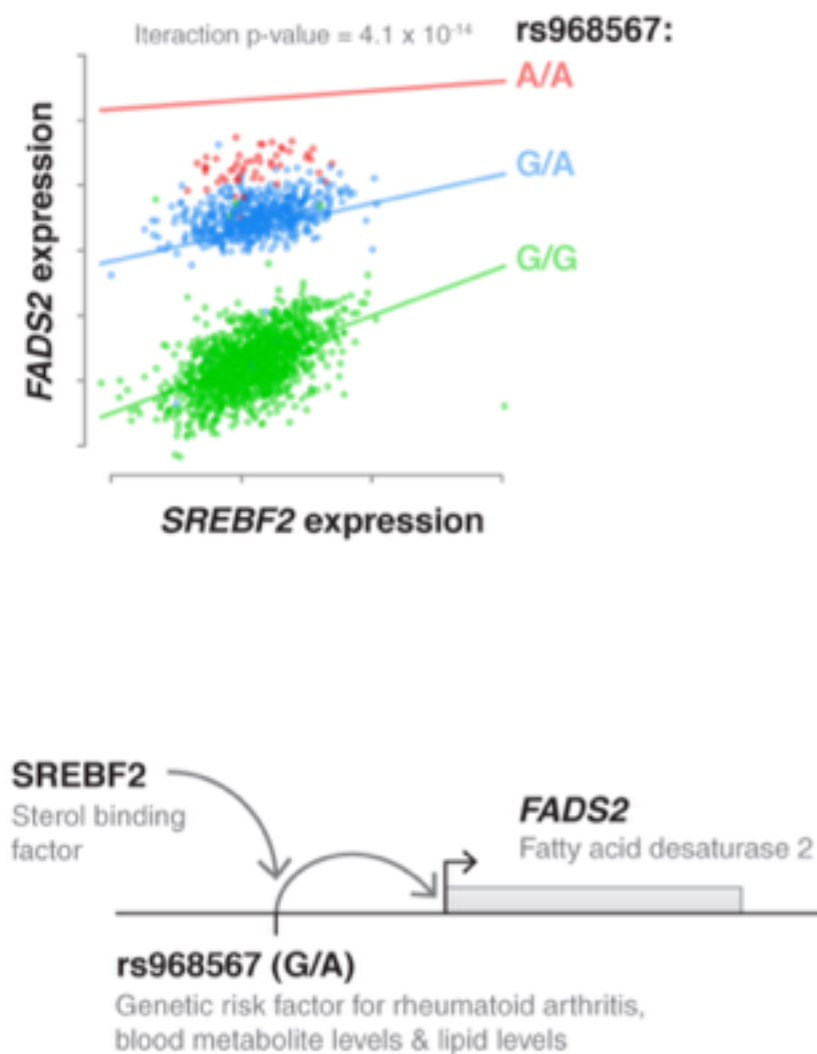low    high

STX3

| Module number | Module description | Number of affected eQTLs | # eQTLs in strong LD with known GWAS hits | GO biological process top enriched pathway |
|---|---|---|---|---|
| 1 | Neutrophils 1 | 917 | 75 | Detection of bacterium |
| 2 | CD4+ T-cells | 337 | 25 | T cell selection |
| 3 | NK cells / CD8+ T-cells | 226 | 19 | Cellular defense response |
| 4 | Erytrocytes | 188 | 8 | Hemoglobin metabolic process |
| 5 | Monocytes / Macrophages | 181 | 11 | Defense response to virus |
| 6 | Growth factor | 156 | 10 | Nerve growth factor receptor signaling pathway |
| 7 | Type 1 interferon | 145 | 11 | Regulation of defense response |
| 8 | Neutrophils 2 | 121 | 3 | Detection of bacterium |
| 9 | B-cells | 123 | 11 | B cell receptor signaling pathway |
| 10 | Eosinophil | 120 | 7 | Regulation of myeloid leukocyte mediated immunity |

Zhernakova *et al*, BiorXiv preprint

Co-expression between top 100 genes per interaction module

Module 7, Top 100 genes

eQTLs with significant interaction with module 7 top covariate gene *SP140*

Zhernakova *et al*, BiorXiv preprint
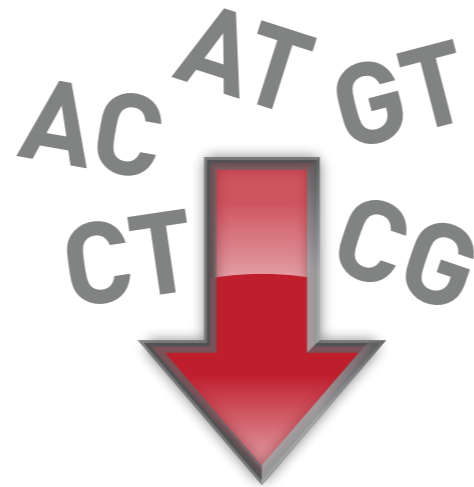
but is this relevant for my patients?

Patient with a severe disease.
You suspect a genetic cause.
What do you do?

- Targeted gene panel?
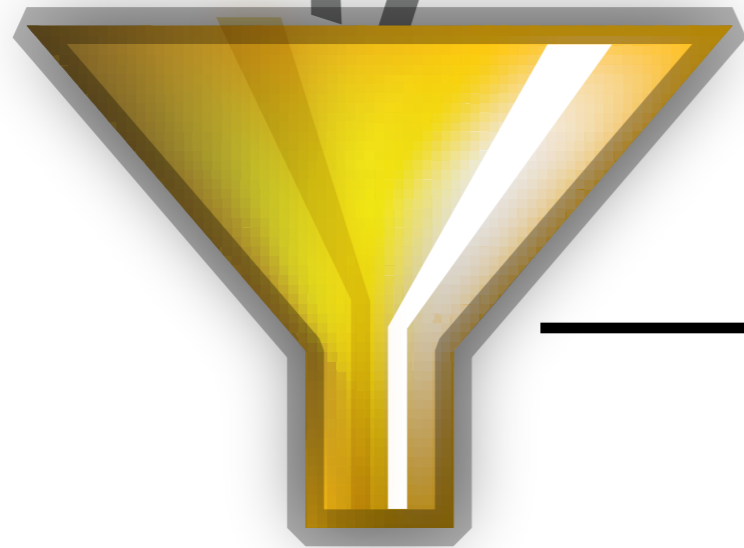- Whole exome sequencing?
- Whole genome sequencing?

Problem:
Many (rare) variants
of unknown significance

_____ gene expression?

- Rare genetic variants also have effects on gene expression

- Rationale BBMRI-NL BIOS Consortium to establish 'Transcriptome of the Netherlands' in 5,000 population based samples

B B M R I • N L

- Generate RNA-seq data on patients. Contrast these expression values to the Transcriptome of the Netherlands.

**TRIM51BP gene expression distribution in the Dutch population**

Log$_2$ expression

Essential to get very accurate reference values for each gene

Get many different cell-types

Recycle big data

Get large sample-sizes

Genes

Samples

Cell-types

Setting:
Activity of switch

Size of switch:
Importance of switch

VOLUME

AMPLIFIER

BALANCE

BASS

MID

TREBLE

L ● ● R

- ● ● +

- ● ● +

- ● ● +

MIN ●

● MAX

PROCESSOR

34SHANNON324

COPYRIGHT 2011

908SGB9477JX

348959AM43

**Wiring:** Way the switch has effect

A control panel that determines gene expression?

Size of switch:
Importance

Setting: State of
a certain sample

Wiring: Effect on
individual genes

TC 4

TC 5

Regulatory factors:
Hormones,
Transcription factors,
Physiological factors,
Other (external) stimuli
Genetic variation

Gene A
Gene B
Gene C
Gene D
Gene E
Gene F
Gene G

Fehrmann *et al*, Nature Genetics 2015

**Component 1**

**Components 1 - 50:**
Physiology, metabolism, cell-type differences

**Component 800**

# GeneNetwork gene function predictions



**GWAS on red blood cell traits:**

Mean hemo-
globin con-
centration:
**rs1175550*G** ▬▬
Chr. 1

*cis*-eQTL
mapping

**Blood eQTL mapping:**

SMIM1 Expression Levels →

P < 10⁻¹⁶

AA    AG    GG

**SMIM1:**
Unknown
function

**Gene function predicton:**
(GeneNetwork.nl, based on
80,000 RNA microarrays)

● Genes known to be involved
   in hemoglobin metabolism

SMIM1

**SMIM1:**
Hemoglobin
metabolism

**Exome sequencing of
individuals, negative
for Vel bloodgroup
antigen:**

AC  AT  GT
CT  CG

Homozygous 17bp
deletion in SMIM1

**Knock-down
in zebrafish:**

Reduced number
of red blood cells

Van der Harst *et al*, Nature 2012                                                    Cvejic *et al*, Nature Genetics 2013

## Amounts of data integrated:

| GWAS in 135,000 samples | eQTL mapping in 1,500 samples | Transcriptomics in 80,000 samples | Exome sequencing | Wet lab proof |

## 697 significant adult height associations:

Wood *et al*, Nature Genetics 2014
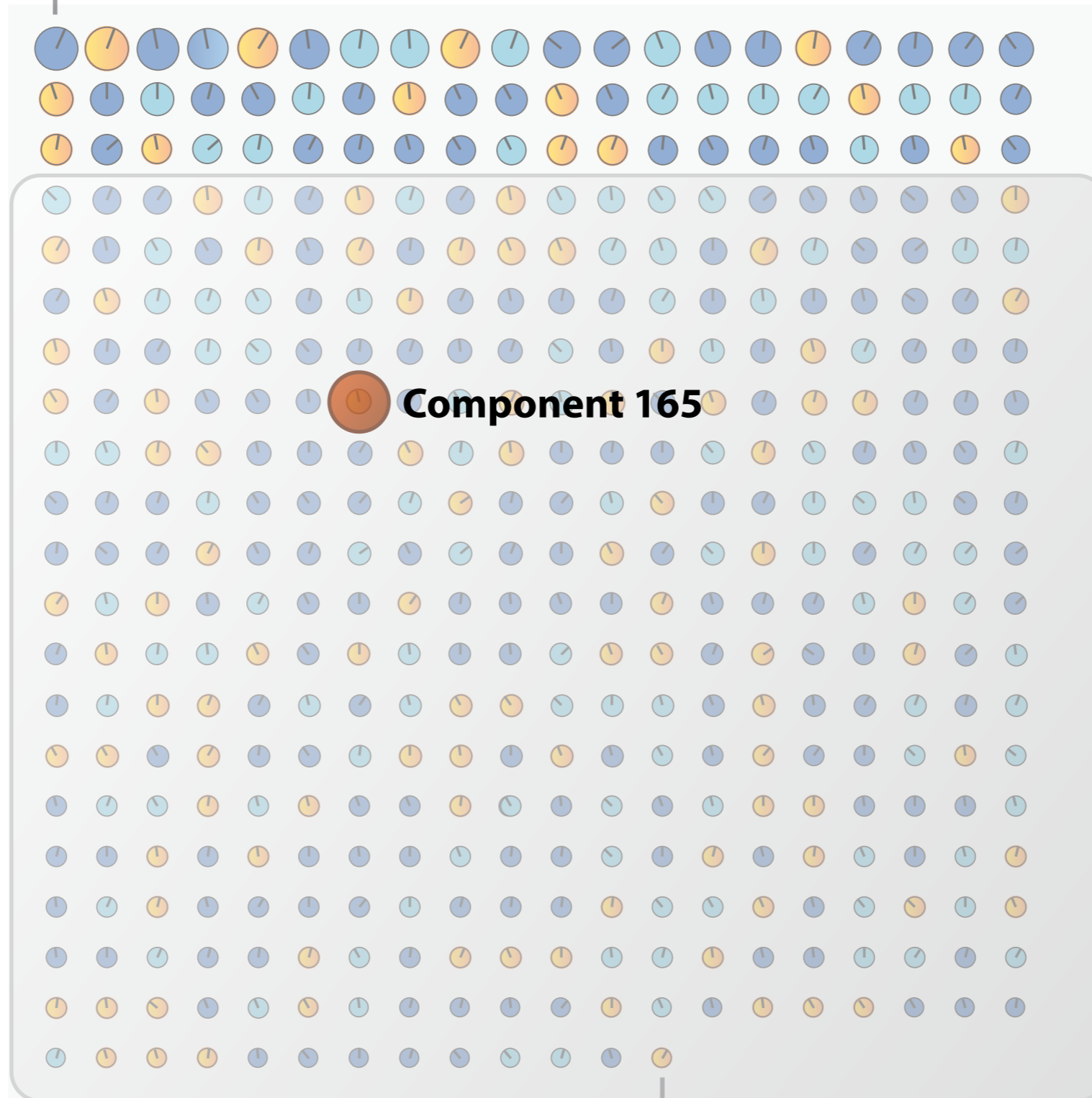


**DEPICT Method:**
Pers *et al*, Nature Communications 2015

**DEPICT used for:**
Body mass index (Locke *et al*, Nature 2015)
Waist hip ratio (Shungin *et al,* Nature 2015)
Hypospadias (Geller *et al*, Nature Genetics 2014)
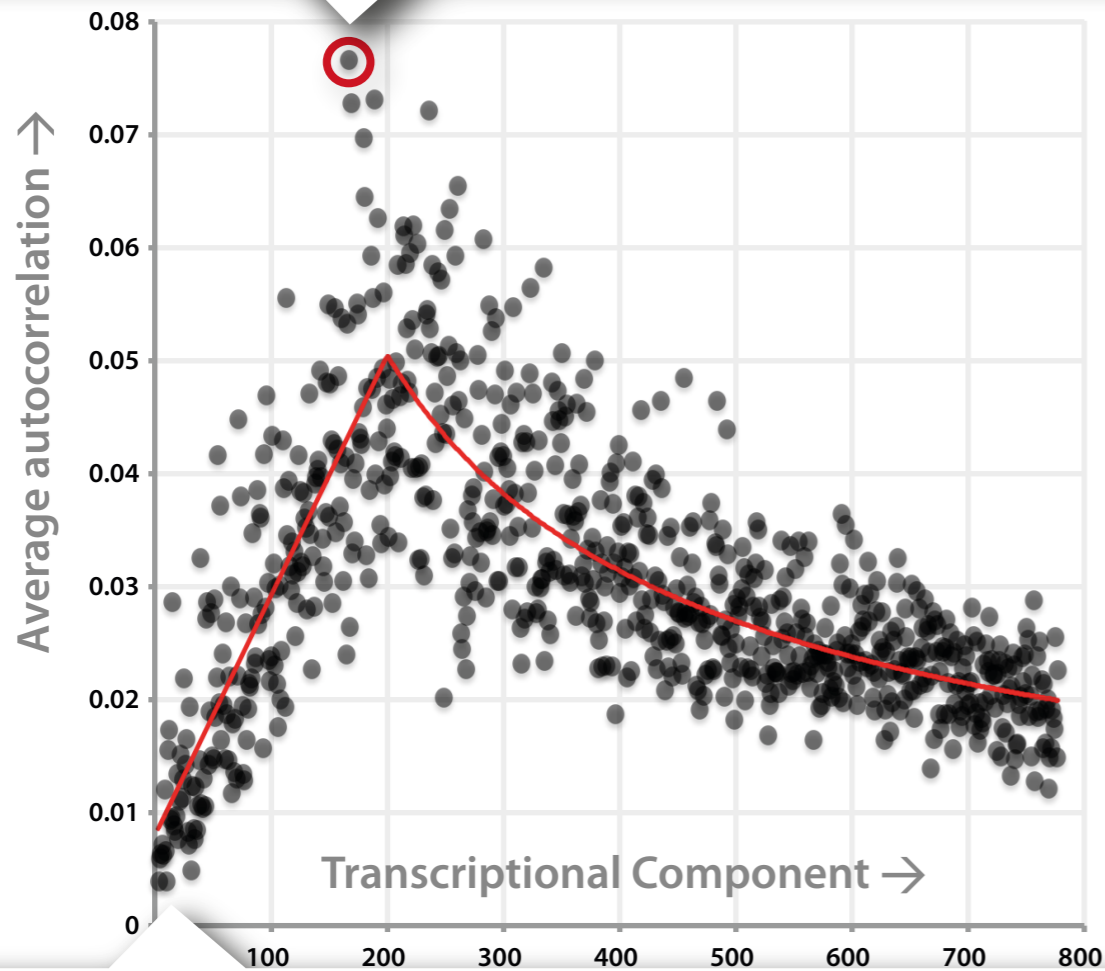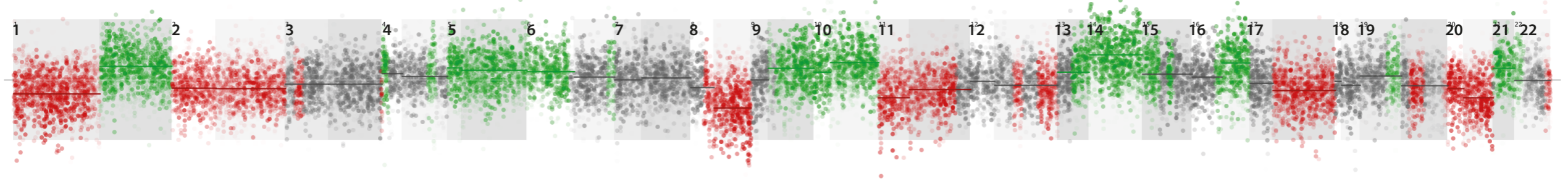Lipid Levels (Surakka, Nature Genetics 2015)

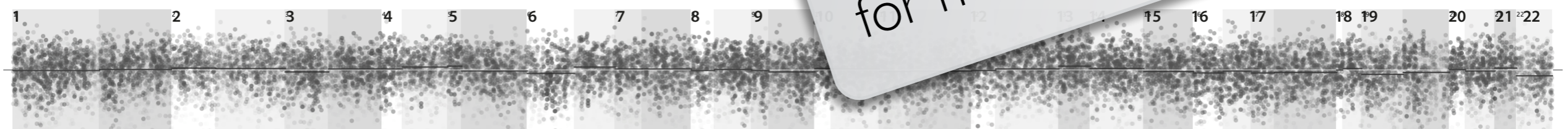Component 1

Component 165

Component 800

# Some component show weird behaviour



**TC 165:** Strong cytogenetic effects, high autocorrelation

**TC 1:** No cytogenetic effect, zero autocorrelation

Redo analysis in healthy samples, correct cancer data for healthy components

Chromosome

Down Syndrome patient: dup 21

Karyogram
HapMap LCL



Chromosome

4   7   9   14   21

GSM274996

cytogenetic RNA expression

arrayCGH

GSM275008

cytogenetic RNA expression

arrayCGH

Fehrmann *et al*, Nature Genetics 2015

Average somatic copy number aberration profile of 16,172 primary tumor samples (GPL570 + GPL96 platforms)

By recycling big data it is possible to clean data and get very accurate measurements

Tipping point at component 165

Average autocorrelation →

Transcriptional Component →

Transition to chaos in the logistic map
Crutchfield *et al*, 1990

Complexity

Entropy

Distribution identical to simulations in complexity theory

**Forest fire:** when will a forest burn down entirely?

How many trees can you plant without the risk that everything burns down?

Percentage of land with living trees after forest fire

100%

50%

0%

0%  50%  100%

Percentage of land filled with trees

20%

Percentage of land with living trees after forest fire

Percentage of land filled with trees
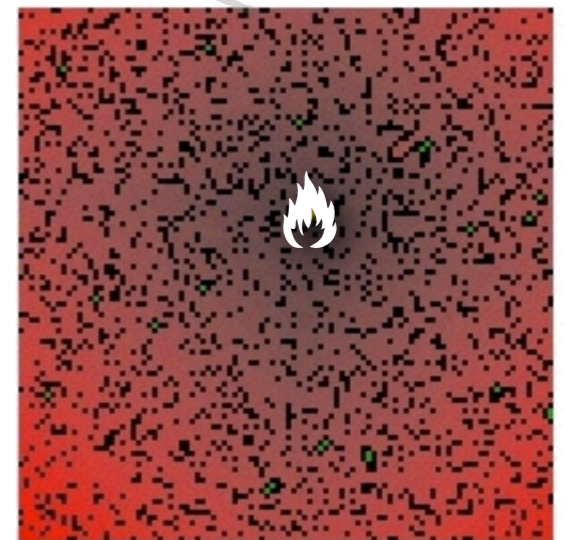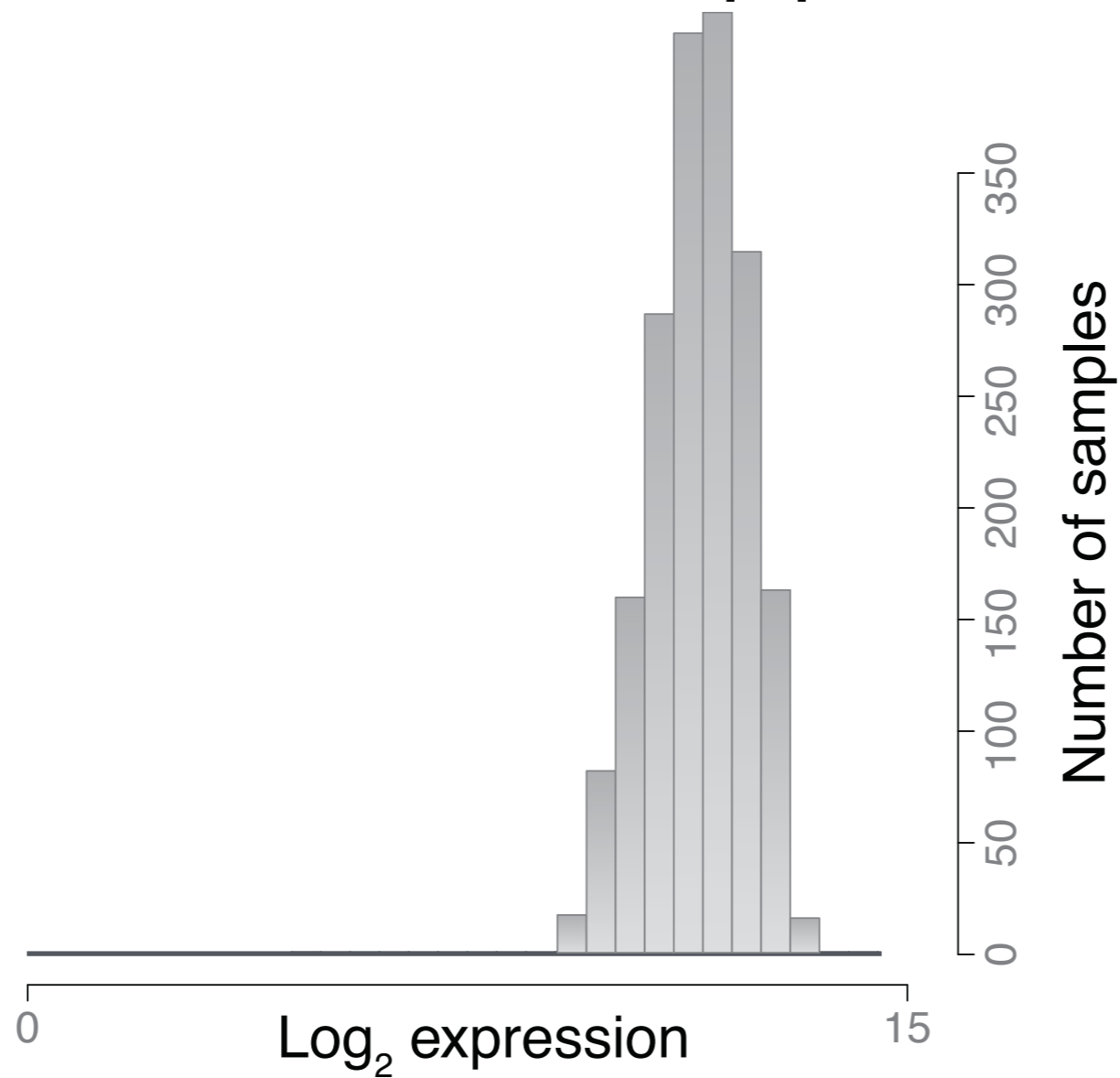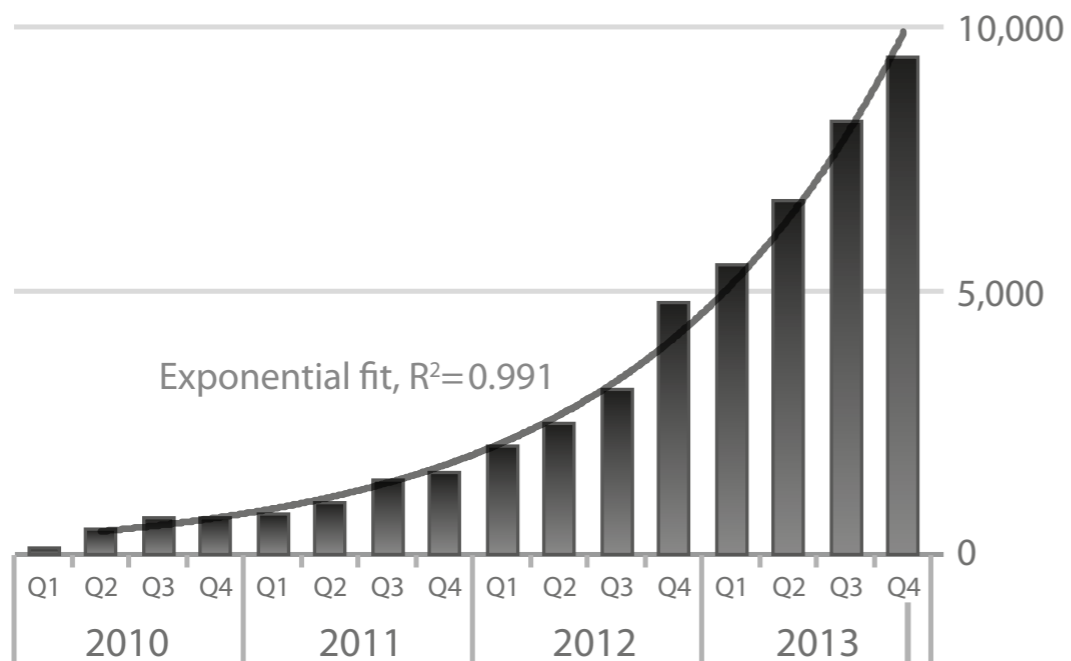
20%

40%

**TRIM51BP gene expression distribution in the Dutch population**

# Explosion of publicly available RNA-seq data

Public RNA-seq data (5,000 samples)
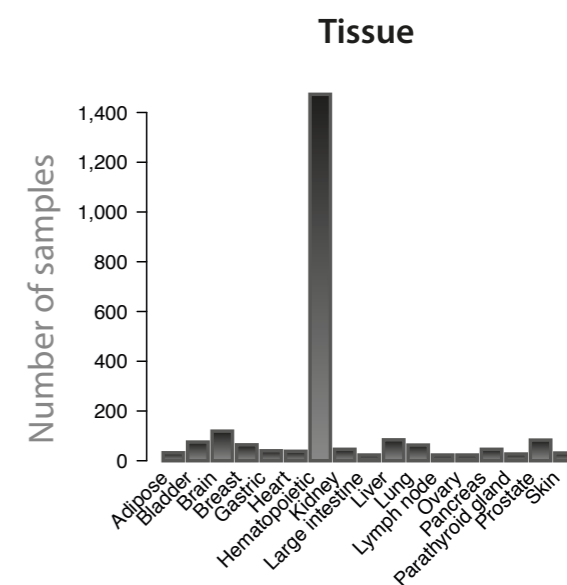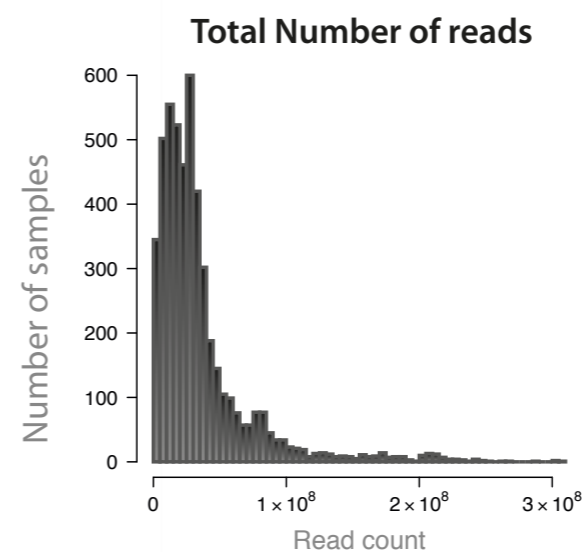
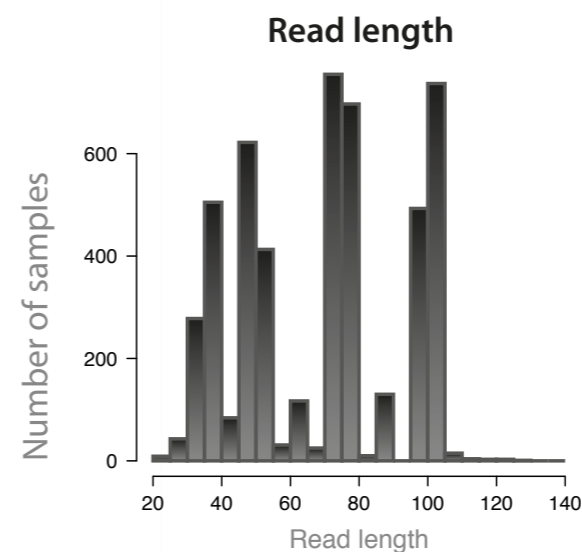Infer genotypes

*KIF13B* expression →

$P = 10^{-21}$

TT    TC    CC

**rs1136055**

Component 2

Component 1

Cell-line

LCLs

B-lympho-cytes

PBMCs

Primary Tissue

Deelen *et al*, Genome Medicine 2015

GATK to call genotypes and output genotype likelihoods, BEAGLE used for imputation towards Genome of the Netherlands

Calling genotypes in RNA-seq data

Ability to call SNP is largely dependent on expressed transcripts

Deelen *et al*, Genome Medicine 2015

Public RNA-seq data: (5,000 samples)

Component 2 / Component 1

Cell-line
LCLs
B-lympho-cytes
PBMCs
Primary Tissue

Deelen et al, Genome Medicine 2015

Genotype calling enables functional effect analysis of:

Common variants: Expression quantitative trait loci

Rare variants: Allele specific expression
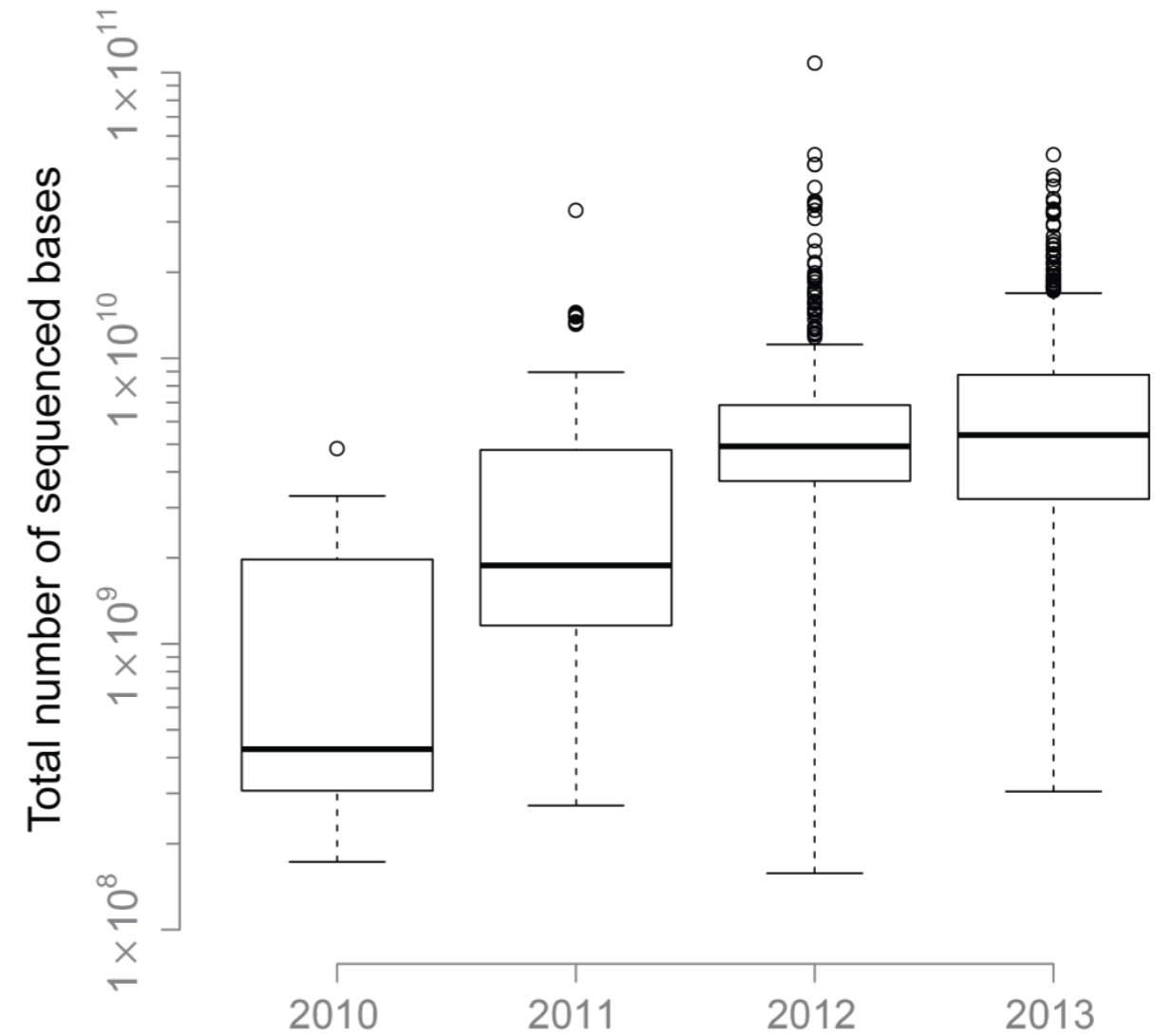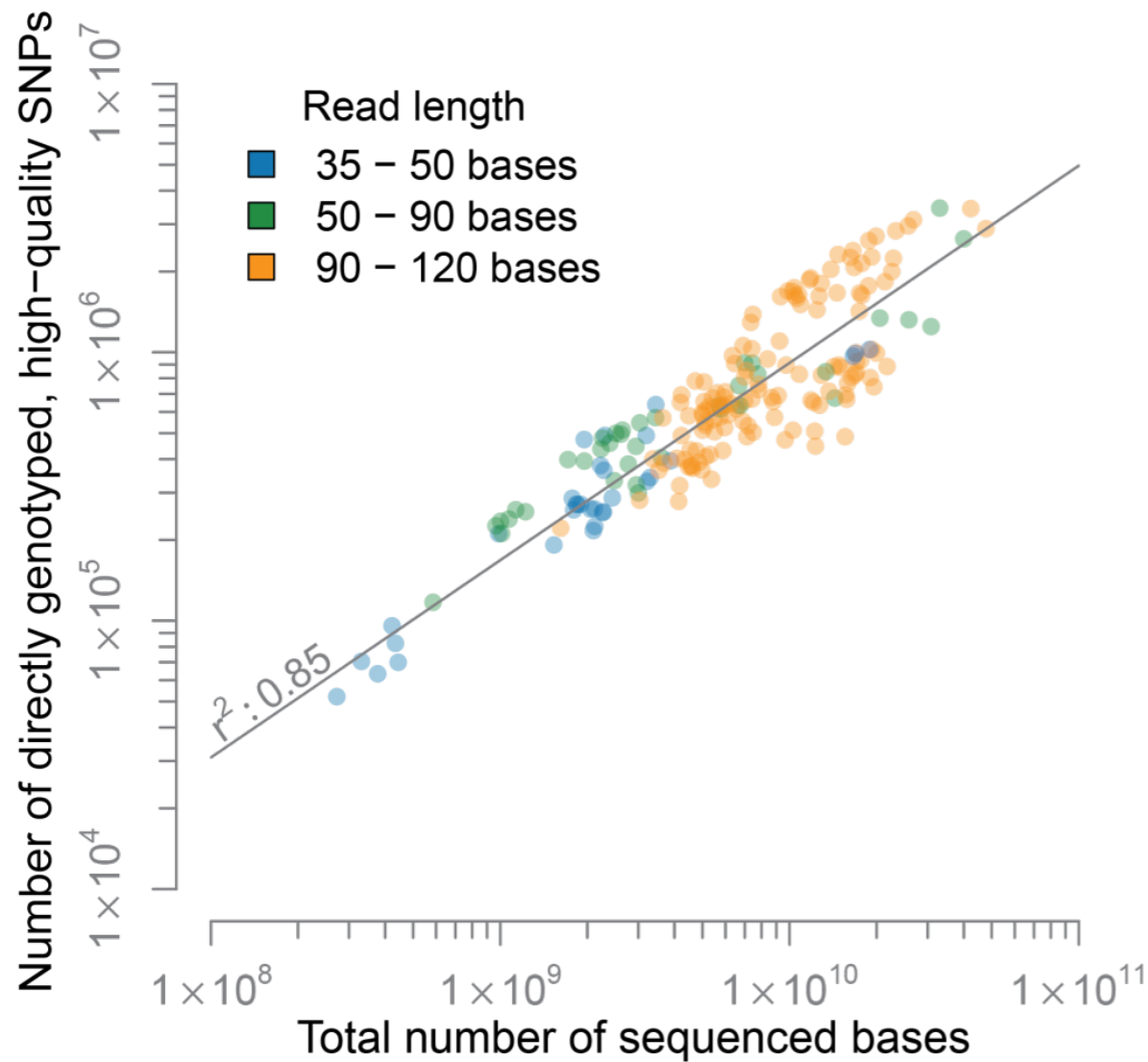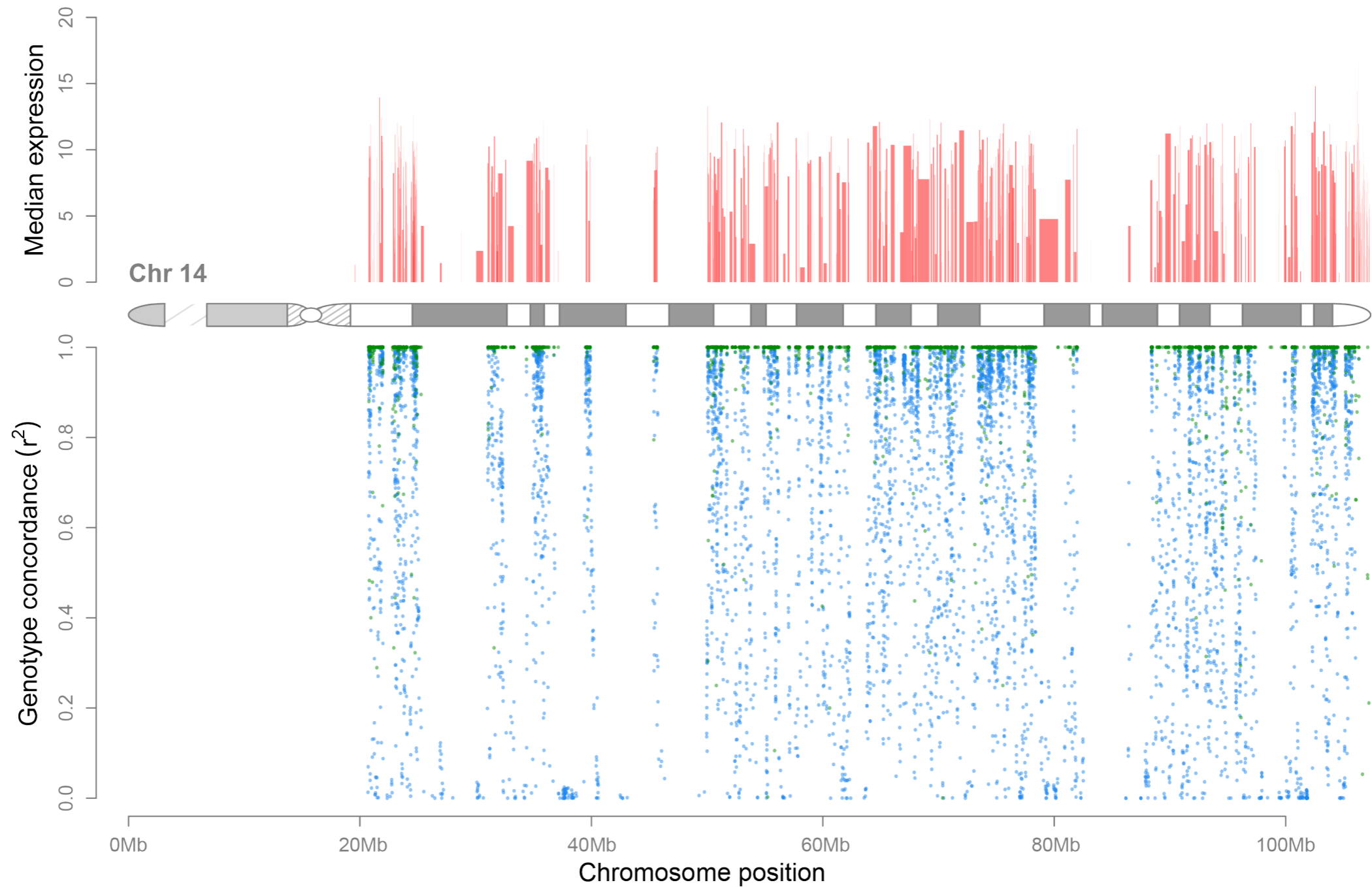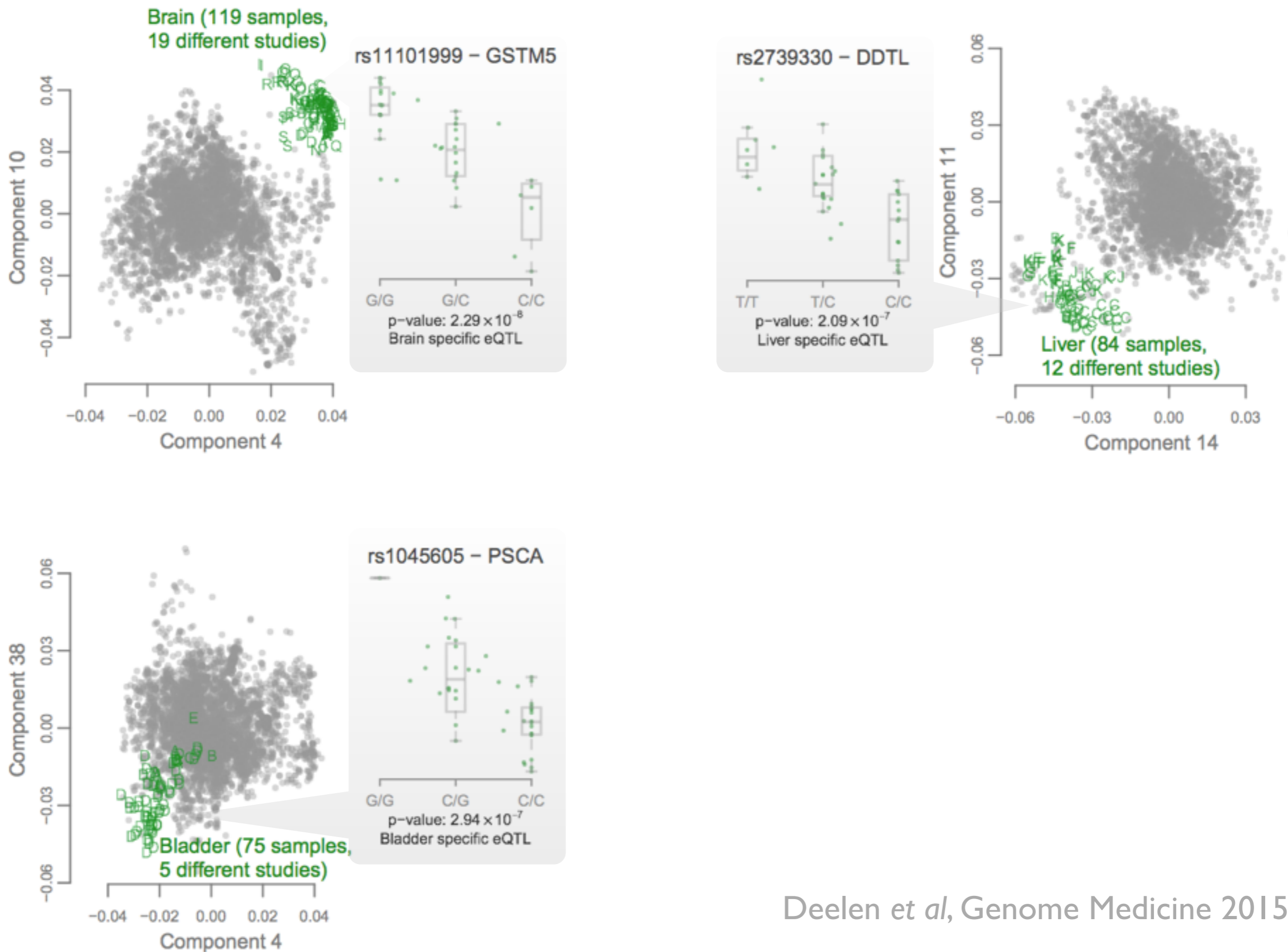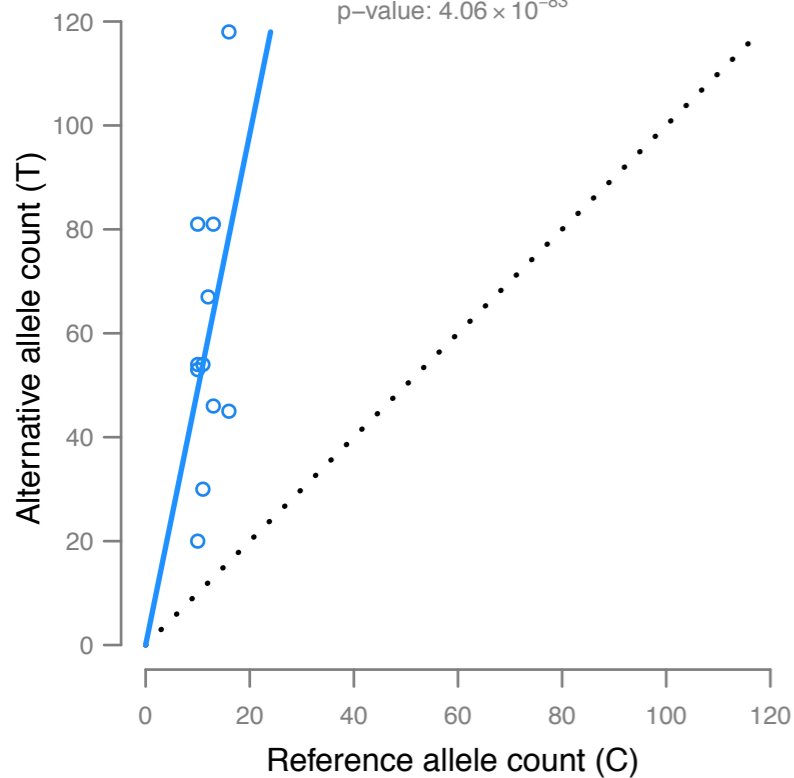
KIF13B expression →

P = 10⁻²¹

TT    TC    CC
rs1136055

rs72550870 – MASP2
p–value: 5.07 × 10⁻¹⁵

Alternative allele count (C)
Reference allele count (T)

Westra et al, Nature Genetics 2013

Uncorrected gene expression profile:

Chromosome  4  7  9  14  21

Gene expression profile, corrected for 'transcriptional components':

Chromosome  4  7  9  14  21

Fehrmann et al, Nature Genetics 2015

Gene expression levels corrected for healthy physiological and metabolic variation

Apply methodology to Individual patients

Apply methodology to Transcriptome of the Netherlands (5,000 samples)

Very low TRIM51BP expression in patient

TRIM51BP gene expression distribution in the Dutch population

Number of samples

Log₂ expression

Candidate causal gene

Patient has certain phenotypes:
- Seizures
- Short stature

Co-regulation identified using public RNA-seq data

TRIM51BP

Genes known to cause seizures

TRIM51BP: co-regulated with known seizure gene

TRIM51BP

Genes known to cause short stature

TRIM51BP: co-regulated with known short stature genes

Candidate causal gene

TRIM51BP likely causal gene
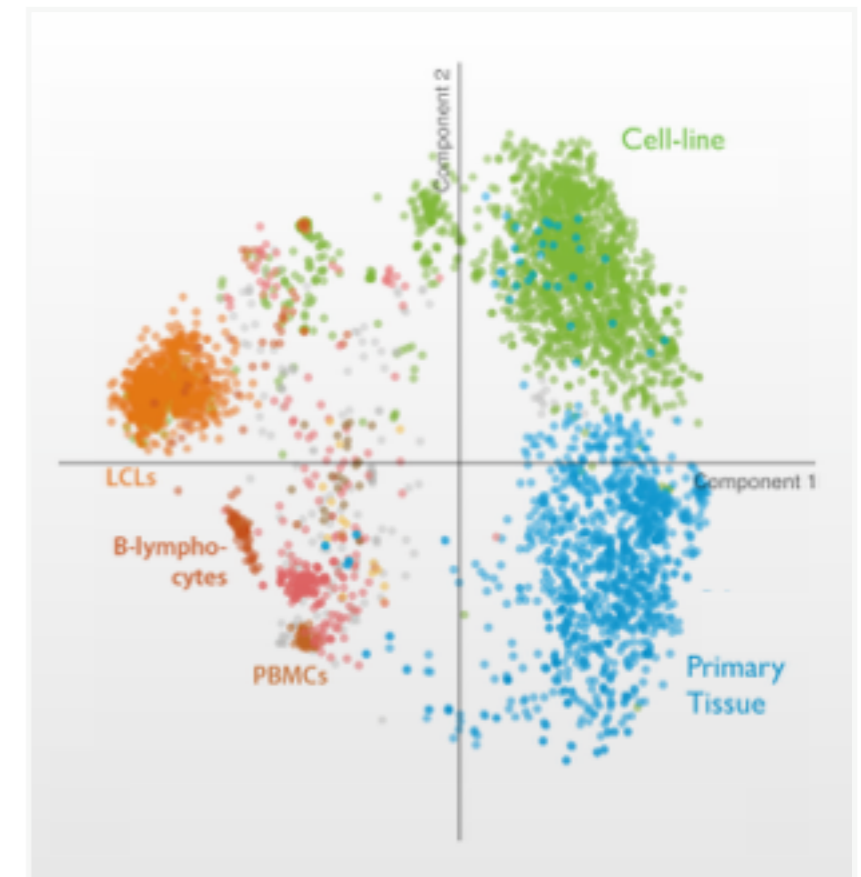
## Lifelines Deep

- 1,500 samples
- Many omics levels
- Genotype data
- Extensive phenotyping

## Transcriptome of the Netherlands

- 5,000 samples
- RNA-seq data
- Genotype data
- Methylation 450k data

## Public RNA-seq data

- 25,000 samples
- RNA-seq data
- Genotype data

- Enormous opportunities exist when recycling 'big data', permits gaining insight into downstream consequences of (rare) genetic variants

- Workshop: how to conduct these analyses yourself:
  - Pointers to the software that is available
  - Identifying sample mix-ups
  - Correcting for unknown confounders
  - Multiple testing correction
  - Allele specific expression

# Acknowledgements >

## Funding >