# Improving the diagnostic yield in clinical genetics by recycling publicly available RNA-seq data

Expression data ——

—— Public data

Phenotypes ——

UMCG
Genetics Department
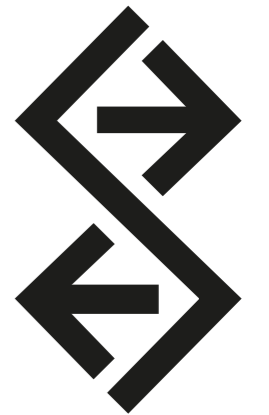
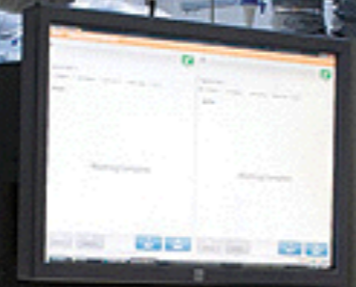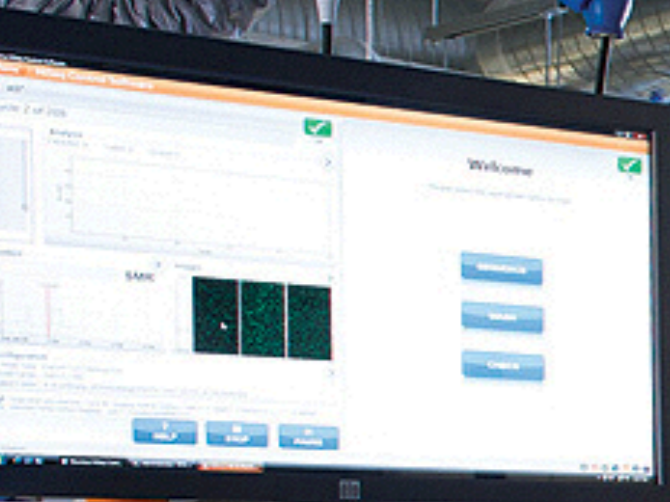'To capture something small you need something big'

CERN

© Ruben van Leer

DNA ——————— AC
CG
GT

'To capture something small
you need something big'

DNA
Sequencers

© Sanger Institute

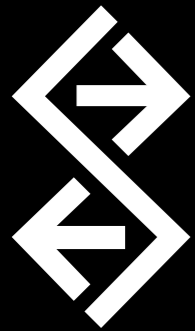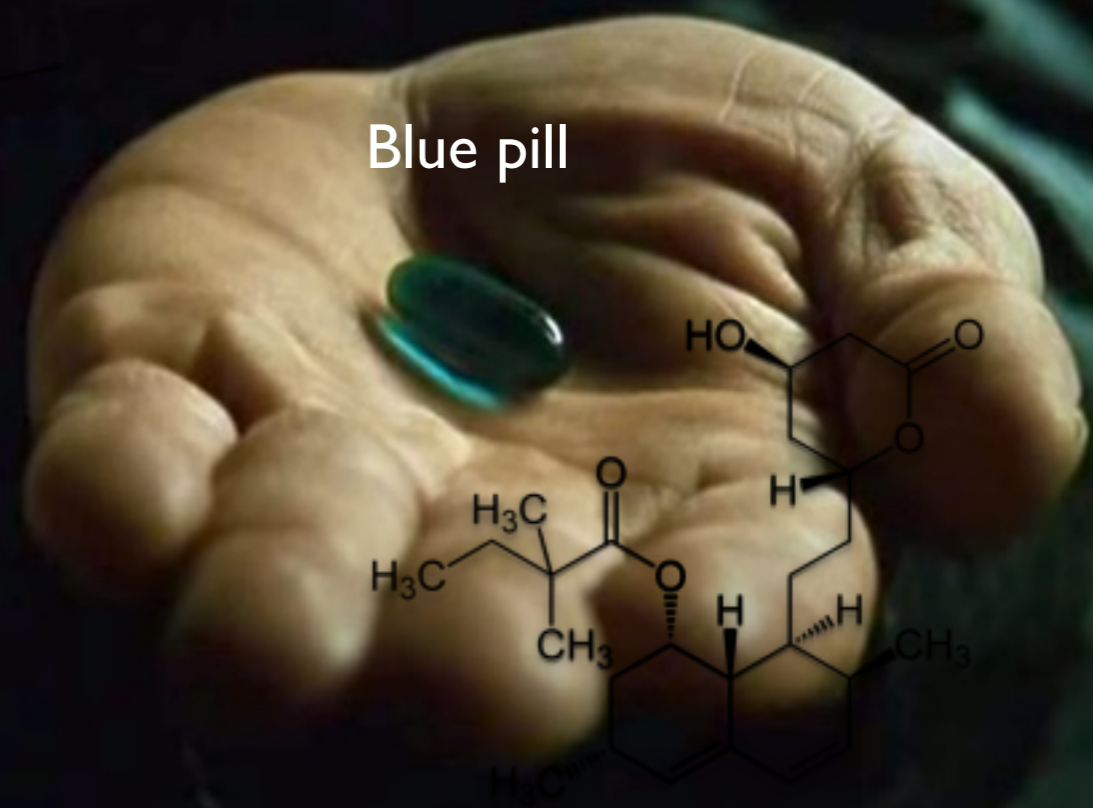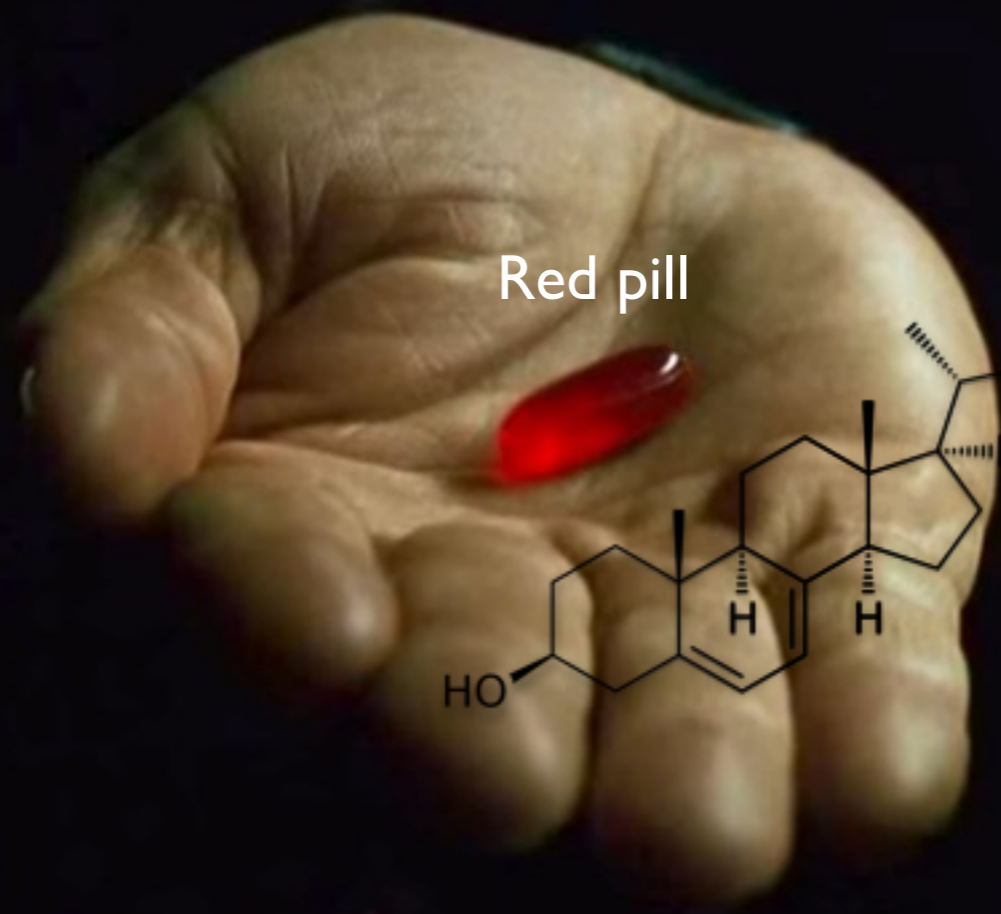'To capture something small you needed something big'

DNA Sequencer

# Minion

© Oxford Nanopore

# 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

Genetic risk factors

Disease

**Black Box**

>10,000 known

Genes unknown
Pathways unknown
Cell-types unknown

>200 diseases

gene expression →

| T-Cell | B-Cell | Monocyte |

P = 0.82

P = 0.42

$P = 10^{-9}$

CC    CT    TT

CC    CT    TT

CC    CT    TT

**genetic risk factor**

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

Genome-wide association studies

cis-eQTL mapping

trans-eQTL mapping

Key driver gene identification

Disease SNP — Disease SNP — Disease SNP — Disease SNP — Disease SNP

cis-eQTL effects:

A    B    C    D    E

trans-eQTL effects:    X    Y

Tissue 1    Tissue 2

Key driver gene    Z

Disease

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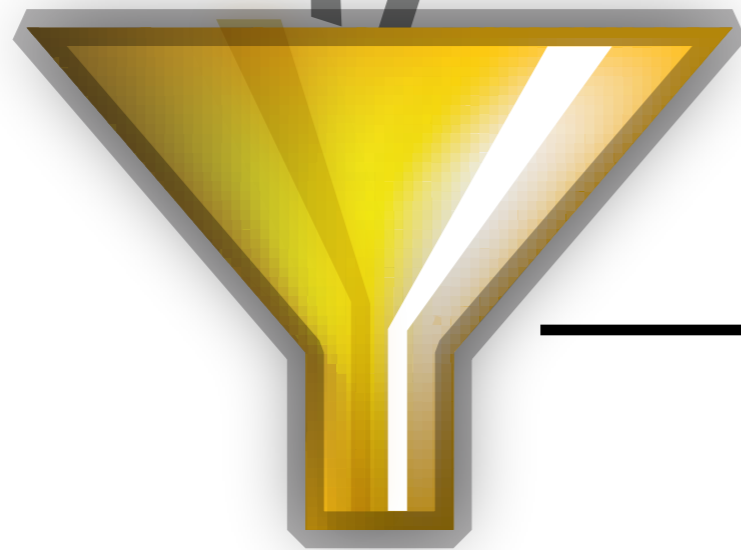
- Rationale BBMRI-NL BIOS Consortium to establish 'Transcriptome of the Netherlands' in 5,000 population based samples

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**

Essential to get very accurate reference values for each gene

## Most expression variation due to:

- Physiological state

- Metabolic state

- Environmental state

RNA blood expression
when you wake up



RNA blood expression
after nice diner

Setting:
Activity of switch

Size of switch:
Importance of switch

VOLUME

AMPLIFIER

BALANCE       BASS       MID       TREBLE

L ●      ● R      - ●      ● +      - ●      ● +      - ●      ● +      MIN ●      ● MAX

PROCESSOR

34SHANNON324       COPYRIGHT 2011

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**

Cell Line
Samples

Blood
Samples

Primary Tissue
Samples

Transcriptional Component 2

Transcriptional Component 1

TC 1    TC 2

# 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)

TC 165: Strong cytogenetic effects, high autocorrelation

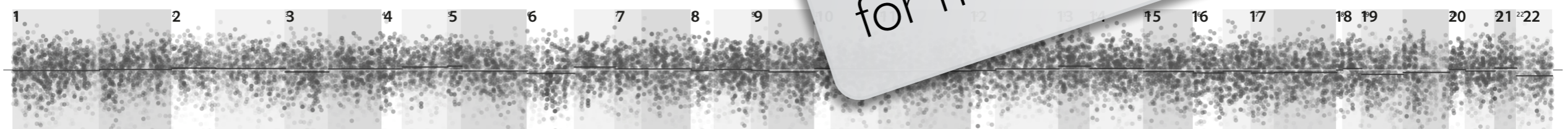TC 1: No cytogenetic effect, zero autocorrelation

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

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

Public RNA-seq data (5,000 samples)

KIF13B expression →

P = 10$^{-21}$

TT    TC    CC

**rs1136055**

Component 2

Cell-line

LCLs

Component 1

B-lympho-
cytes

Primary
Tissue

PBMCs

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

Component 2

Cell-line

Component 1

LCLs

B-lympho-
cytes

PBMCs

Primary
Tissue

Deelen et al,
Genome Medicine 2015

Genotype calling enables
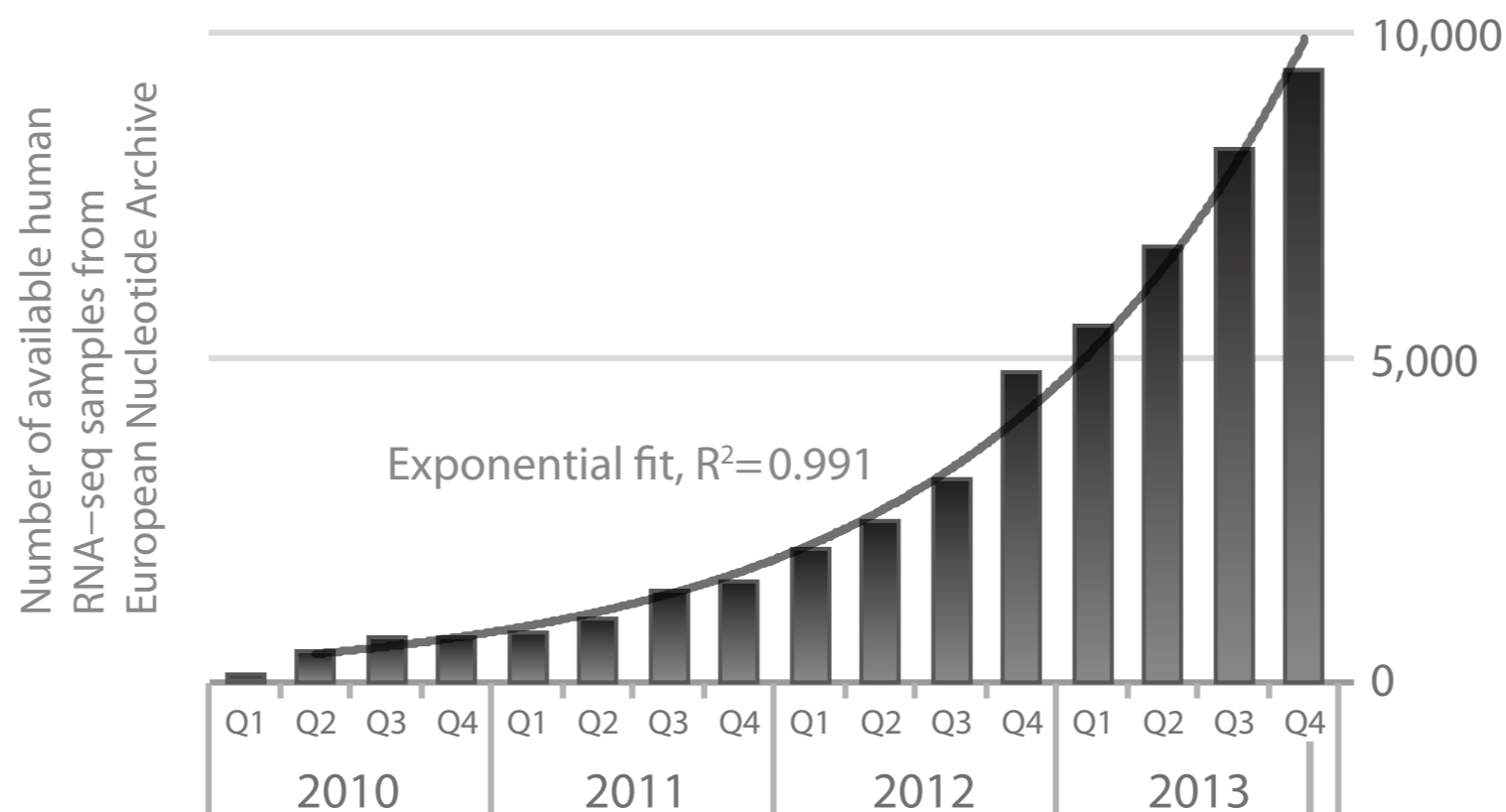functional effect analysis of:

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)

400

300

200

100

0

0        100        200        300        400
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

350
300
250
200
150
100
50

0        Log$_2$ expression        15

**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**

# Acknowledgements >

## Funding >