

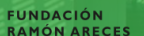
# Introduction to omic data

Juan R González  
Bioinformatics Research Group in Epidemiology

**ISGlobal**  
Institut de  
Salut Global  
Barcelona



Una iniciativa de:



## OUTLINE

- ☐ Omic data
- ☐ Genomics, transcriptomics and epigenomics
- ☐ Omic databases
- ☐ Bioconductor



## Omic data

### OMIC

- ❑ The term “omic” is derived from the Latin suffix “ome” meaning mass or many. Thus, OMICS involve a mass (large number) of measurements/features per outcome (Jackson et al., 2006)

### Integration of OMICS data

- ❑ Efficient integration of data from different OMICS can greatly facilitate the discovery of true causes and states of disease, mostly done by softwares (Andrew et al., 2006).

## Omic data

- ❑ In biological context , suffix **–omics** is used to refer to the study of **large sets of biological molecules** (Smith et al., 2005)
- ❑ The realization that DNA is not alone regulate complex biological processes (as a result of HGP, 2001), triggered the rapid development of **several fields in molecular biology** that together are described with the term OMICS.
- ❑ The OMICS field ranges from
  - ❑ Genomics (focused on the genome)
  - ❑ Proteomics (focused on large sets of proteins, the proteome)
  - ❑ Metabolomics (focused on large sets of small molecules, the metabolome).

## Omic data

The field of genomics has been divided into 3 major categories:

- ❑ **Genotyping** (focused on the genome sequence): the physiological function of genes and the elucidation of the role of specific genes in disease susceptibility (Syvanen, 2001) [Structural genomics]
- ❑ **Transcriptomics** (focused on genomic expression): The abundance of specific mRNA transcripts in a biological sample is a reflection of the expression levels of the corresponding genes (Manning et al., 2007) [Functional genomics]
- ❑ **Epigenomics** (focused on epigenetic regulation of genome expression): Study of epigenetic processes (expression activities not involving DNA) on a large (ultimately genome-wide) scale (Feinberg, 2007) [Functional genomics]

# Genotyping

## Goal

- ❑ Identification of the physiological function of genes
- ❑ Role of specific genes in disease susceptibility (syvanen et al., 2001)

## Common feature used

- ❑ Among different variations (insertions, deletions, SNPs, etc.), **single nucleotide polymorphisms** (SNPs) are the most commonly investigated (Sachidanandam et al., 2001) and can be used as markers for diseases.
- ❑ Tag SNPs (informative subset of SNPs) and fine mapping are further used to identify true cause of phenotype (patil et al., 2001).

## Application

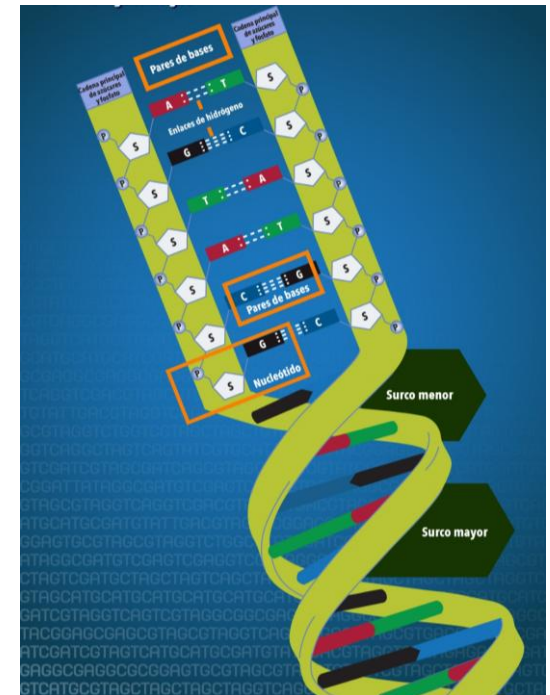
- ❑ Identification of genes associated with disease

## Recent improvement in genotyping

- ❑ **Array-based genotyping** allows the assessment of entire genome (up to 1M SNPs) per assay: GWAS (Jelly et al., 2010)
- ❑ Next Generation Sequencing (**WES, WGS**)

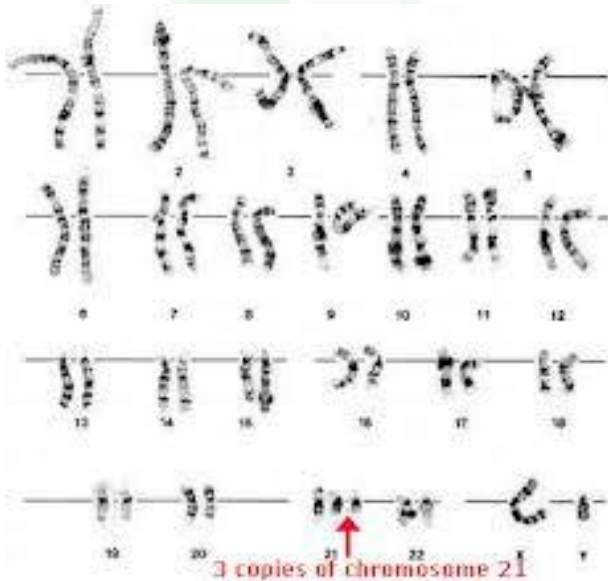
# Genotyping

- ❑ DNA = deoxyribonucleic acid
- ❑ Two antiparallellele strands: double helix
- ❑ Four nucleotides:
  - Adenine = Timine
  - Cytosine = Guanine
- ❑ Human genome size: 3.2 billion bp
- ❑ Known variants: 324 M variants (2017)
- ❑ Difference between 2 individuals: 20 M bp (0.6%)

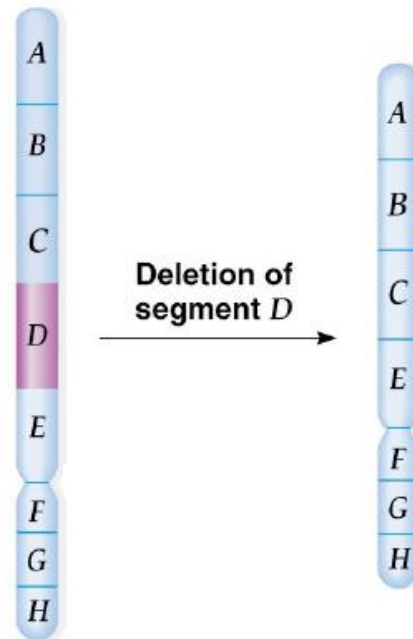


# Genotyping

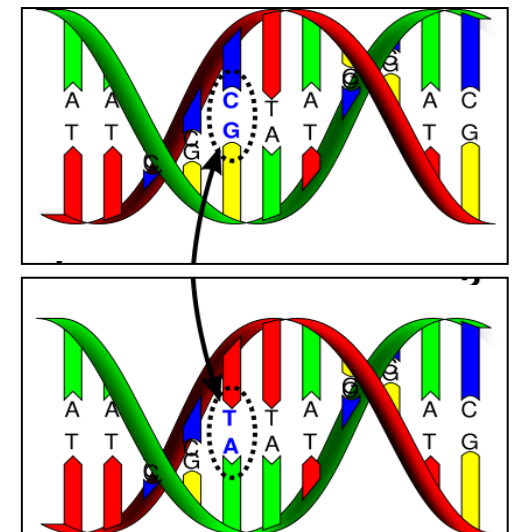
## Chromosomal variants



## Structural variants

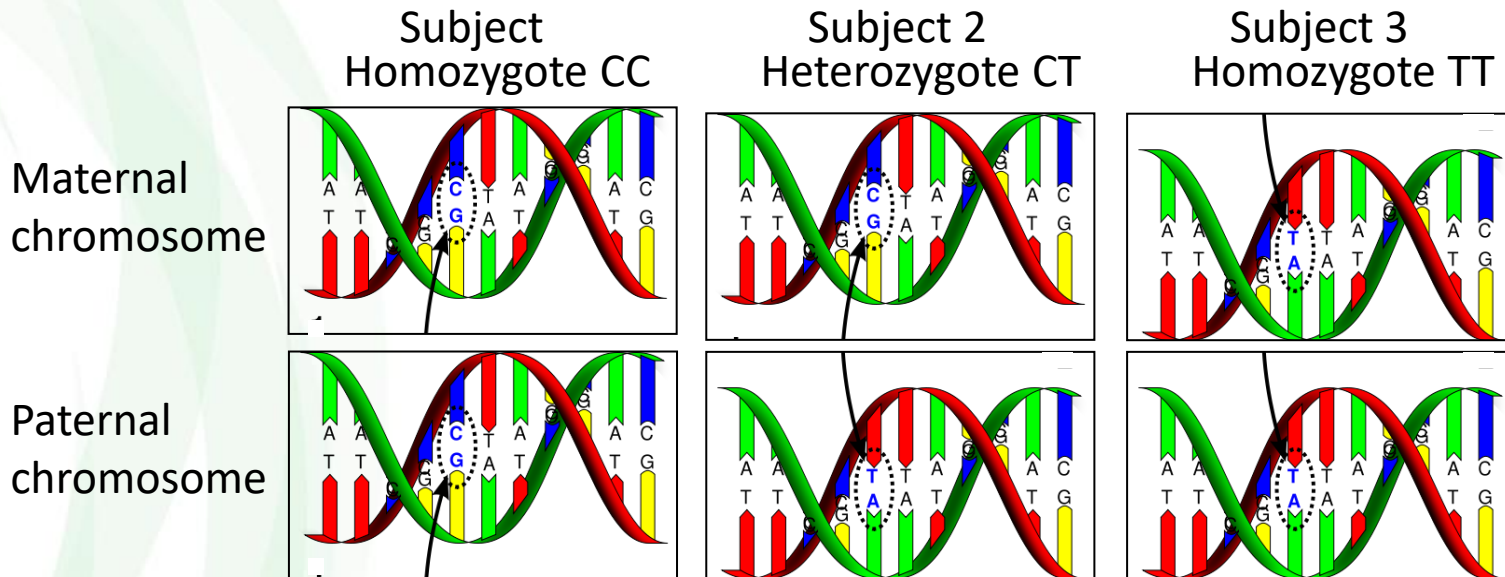


## Single nucleotide polymorphism (SNP)





# Genotyping



Genotype vs allele

Allele: different forms a genetic variant can take

Locus/Loci

# Genotyping

## Human Reference Genome

- A single consensus representation of the genome

<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>

## Assembly (versions):

- Human genome 19 (hg19)=build 37 (bd37)=GRCh37
- Human genome 38 (hg38)=GRCh38

## Strand:

- genome assembly strand (plus (+) or minus (-))
- dbSNP strand (forward (F) or reverse (R))
- Illumina (top (T) or bottom (B))

# Genotyping

## Human genome browsers

- System for displaying, viewing and accessing genome annotation data
- NCBI (National Center for Biotechnology Information):  
<http://www.ncbi.nlm.nih.gov/>
- ENSEMBL (a joint project between EMBL (European Molecular Biology Laboratory)-EBI and the Wellcome Trust Sanger Institute):  
<http://www.ensembl.org/index.html>
- UCSC (University of California Santa Cruz) Genome Bioinformatics:  
<http://genome.ucsc.edu/>

# Genotyping

## 1000 Genomes Project

### ARTICLE

---

---

doi:10.1038/nature09534

## A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium\*

### ARTICLE

---

---

doi:10.1038/nature11632

## An integrated map of genetic variation from 1,092 human genomes

The 1000 Genomes Project Consortium\*



# Genotyping

## Genetic variants database: SNPs

- dbSNP: <http://www.ncbi.nlm.nih.gov/projects/SNP>
- ENSEMBL:  
[http://www.ensembl.org/info/website/tutorials/gene\\_snps.html](http://www.ensembl.org/info/website/tutorials/gene_snps.html)

SNPs have a unique identifier “rs#”

Ex. rs1695

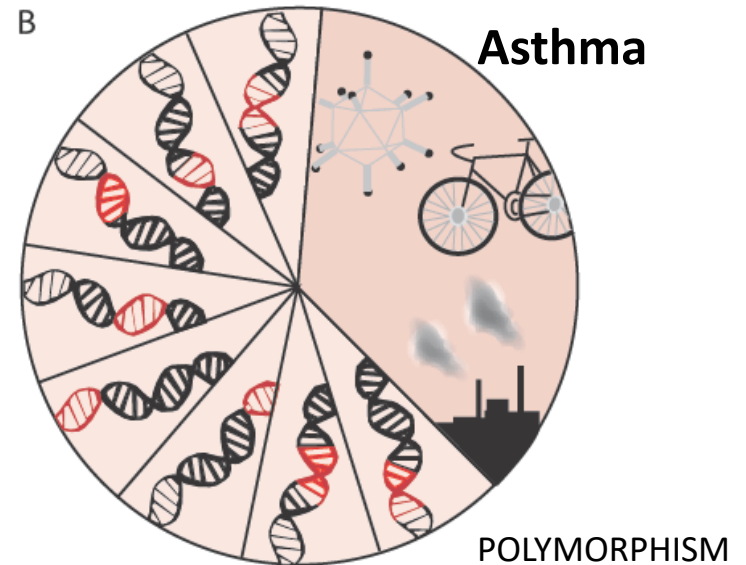
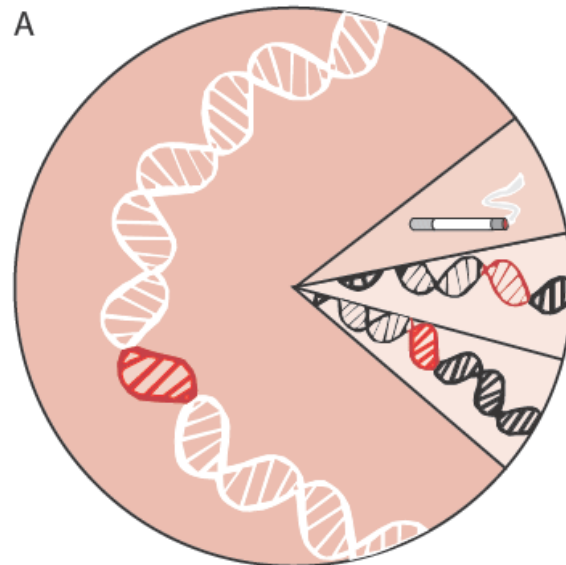
# Genotyping

## Mendelian vs complex diseases

### Phenylketonuria

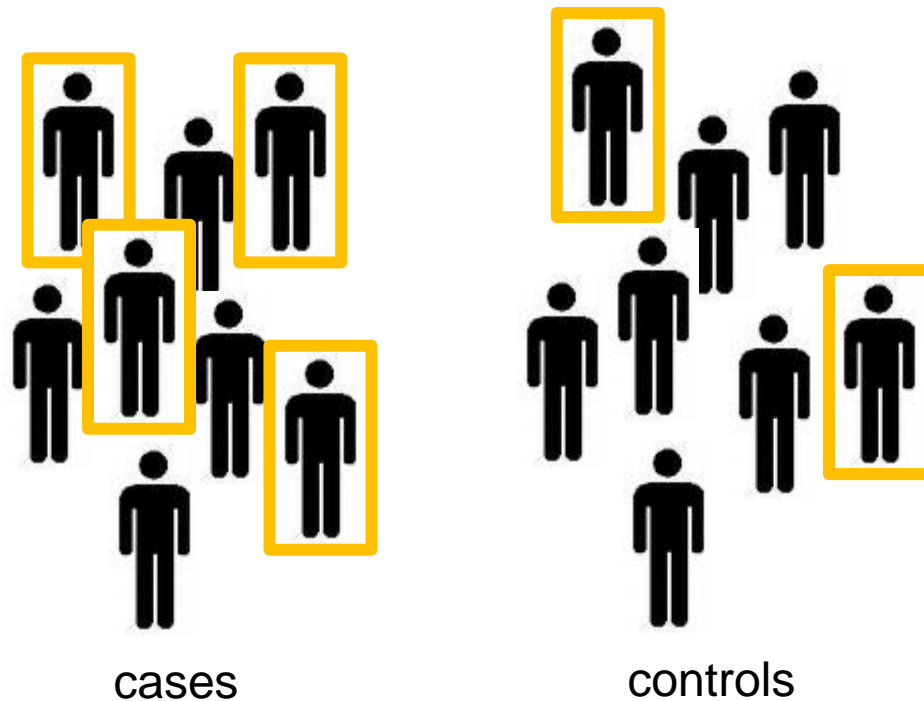
- mutations in PAH gene
- decreased metabolism of phenylalanine
- treatment: food with low levels of phenylalanine

MUTATION



# Genotyping

**Association studies:** Genetic marker vs. phenotypic trait



Variant (SNP)  
allele

# Genotyping

- ☐ Hypothesis driven study: SNP, gene/locus, pathway...
- ☐ Agnostic screening
- ☐ Genome wide association study (GWAS) – arrays
- ☐ Whole genome/exome sequencing (WGS, WES)

## GENOTYPING

Particular points in the genome

....AGCTAAATGATAGCATCAT....

## SEQUENCING

All the nucleotides in a genomic region

....AGCTAAATGATAGCATCAT....



# Genotyping

GWAS catalog: <https://www.ebi.ac.uk/gwas/>

## REVIEW

### 10 Years of GWAS Discovery: Biology, Function, and Translation

Peter M. Visscher,<sup>1,2,\*</sup> Naomi R. Wray,<sup>1,2</sup> Qian Zhang,<sup>1</sup> Pamela Sklar,<sup>3</sup> Mark I. McCarthy,<sup>4,5,6</sup>  
Matthew A. Brown,<sup>7</sup> and Jian Yang<sup>1,2</sup>

Application of the experimental design of genome-wide association studies (GWASs) is now 10 years old (young), and here we review the remarkable range of discoveries it has facilitated in population and complex-trait genetics, the biology of diseases, and translation toward new therapeutics. We predict the likely discoveries in the next 10 years, when GWASs will be based on millions of samples with array data imputed to a large fully sequenced reference panel and on hundreds of thousands of samples with whole-genome sequencing data.

# Genotyping

## ❑ Selection of individuals

❑ Phenotype

❑ DNA quality

❑ **Population stratification** (ethnicity, technical...)

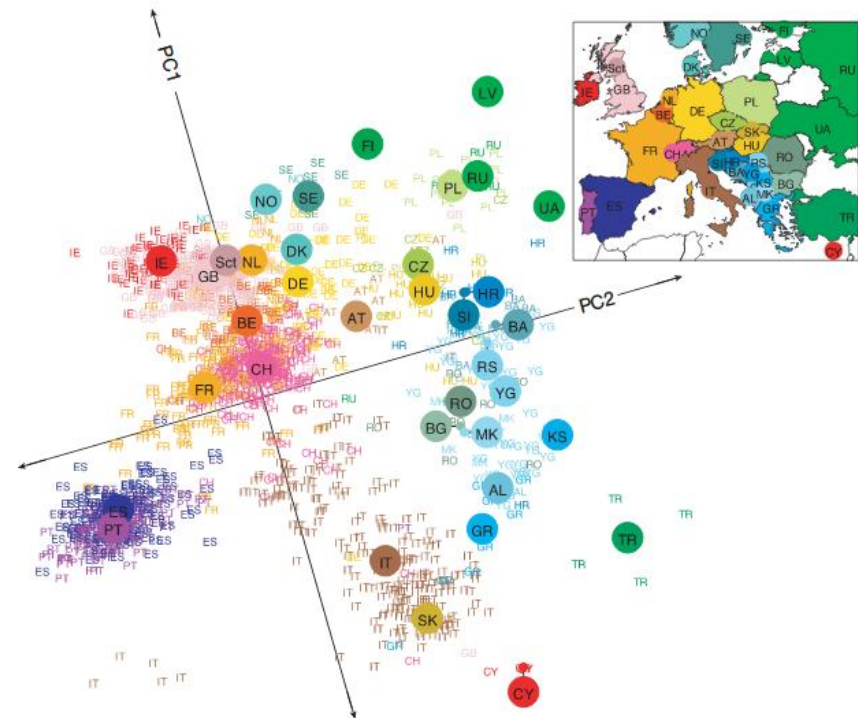
Ethnicity:

European and African different genetic background

European (in blue) high prevalence of disease X

African low prevalence of disease X

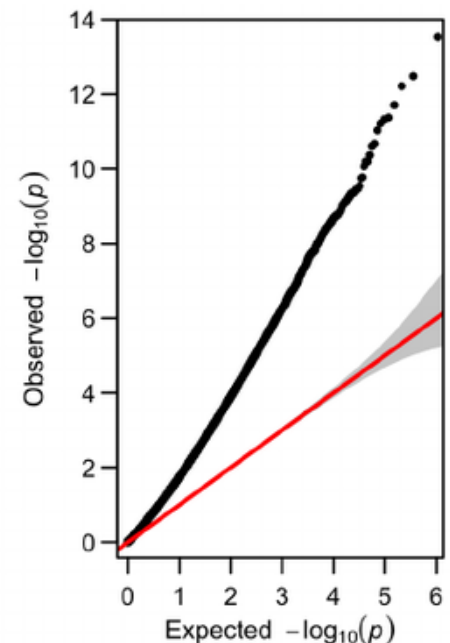
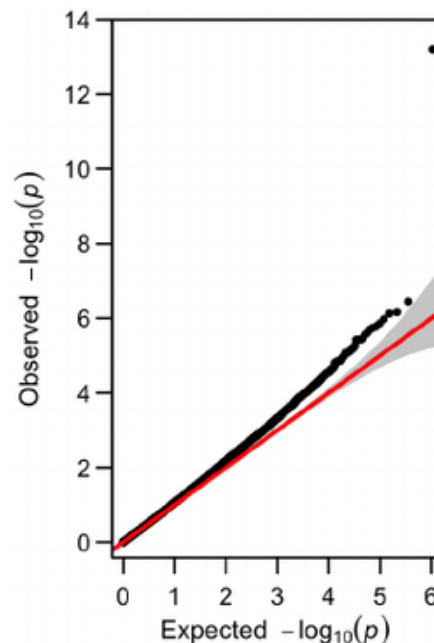
Novembre et al 2008



# Genotyping

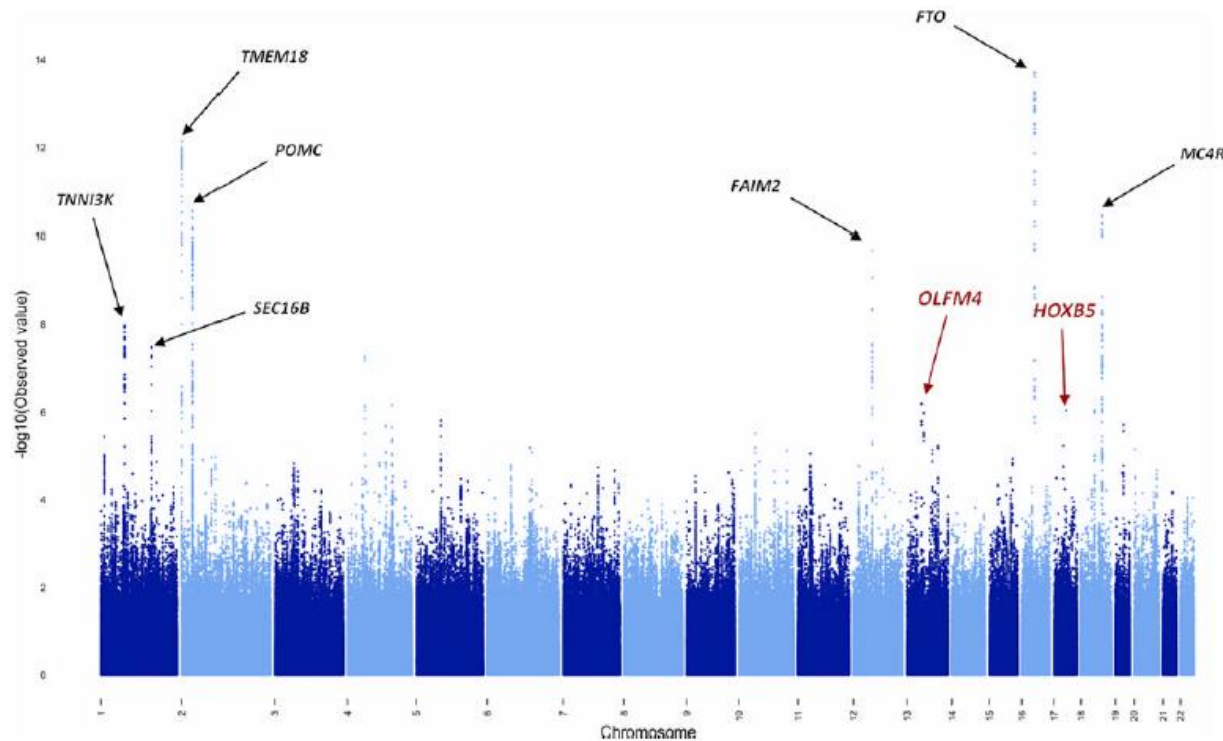
## GWAS multipletesting correction

- ❑ Many variants, then multipletesting correction needed
- ❑ Genome wide significant:  
P value < 5E-08 (number of independent LD blocks in genome)
- ❑ Suggestive association:  
P value < 1E-05



# Genotyping

## Manhattan plot



**Figure 1.** Manhattan Plot of the meta-analysis of childhood obesity GWAS runs in the discovery stage (5,530 cases and 8,318 controls), with each locus achieving genome wide significance ( $P < 5 \times 10^{-8}$ ) indicated in black text. In addition, the novel loci uncovered in this study are indicated in red text.



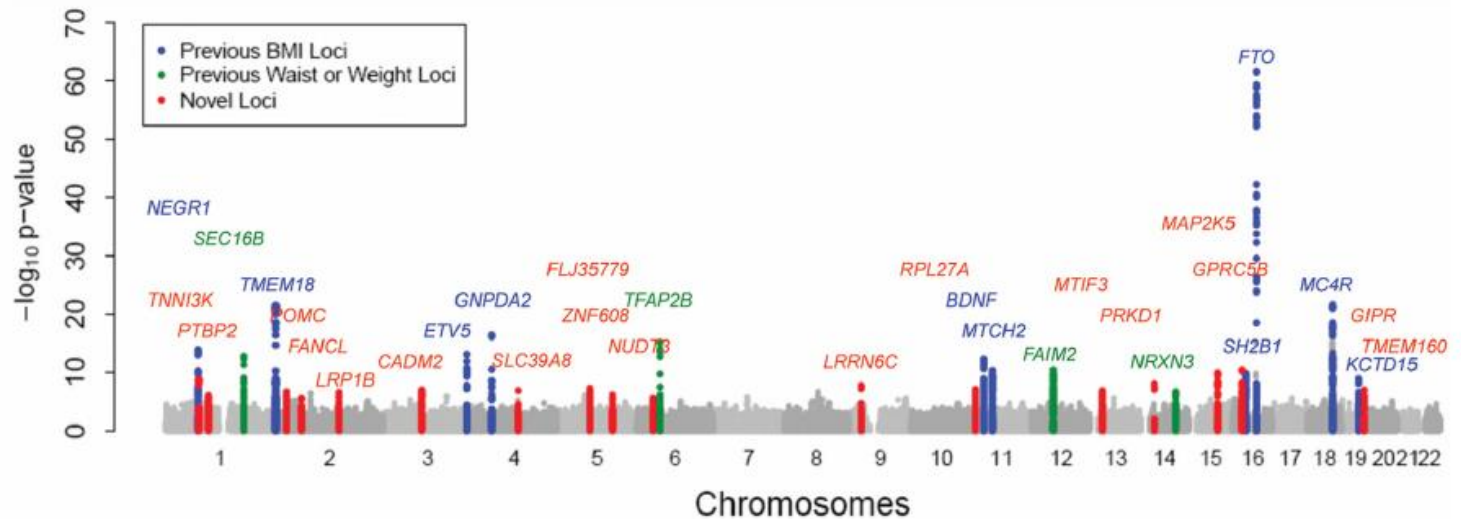
# Genotyping

## Personalized medicine

*Nat Genet.* 2010 November ; 42(11): 937–948. doi:10.1038/ng.686.

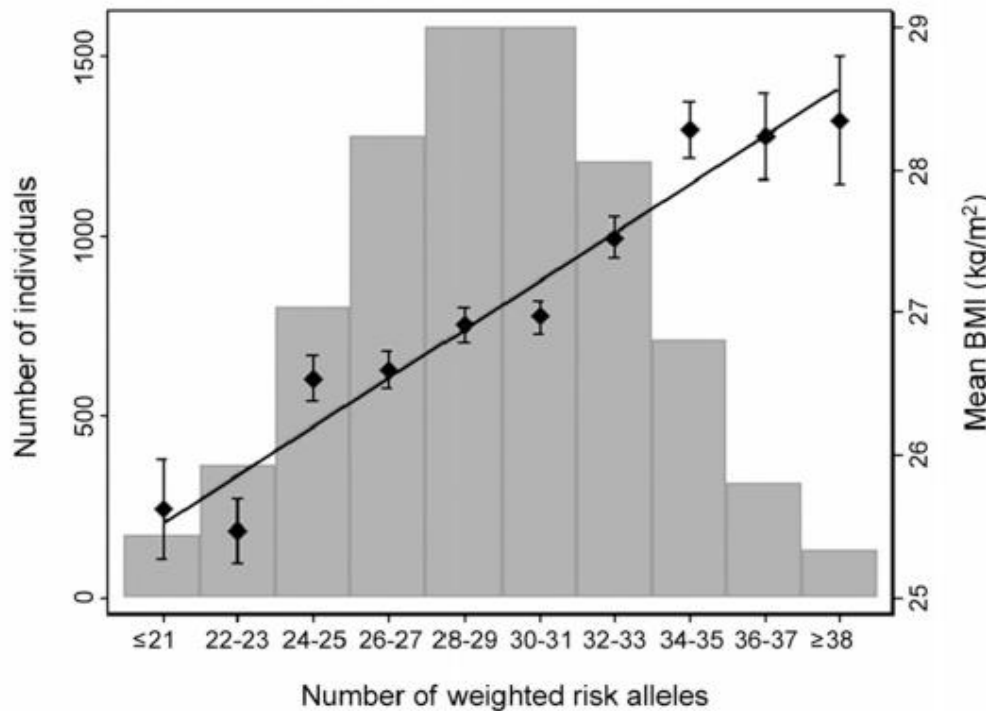
### Association analyses of 249,796 individuals reveal eighteen new loci associated with body mass index

Elizabeth K. Speliotes<sup>1,2,\*</sup>, Cristen J. Willer<sup>3,\*</sup>, Sonja I. Berndt<sup>4,\*</sup>, Keri L. Monda<sup>5,\*</sup>, Gudmar Thorleifsson<sup>6,\*</sup>, Anne U. Jackson<sup>3</sup>, Hana Lango Allen<sup>7</sup>, Cecilia M. Lindgren<sup>8,9</sup>, Jian'an Luan<sup>10</sup>, Reedik Mägi<sup>8</sup>, Joshua C. Randall<sup>8</sup>, Sailaja Vedantam<sup>1,11</sup>, Thomas W. Winkler<sup>12</sup>, Lu



# Genotyping

## Personalized medicine



ROC curve

# Transcriptomics

## Gene expression profiling

- ❑ The identification and characterization of the mixture of mRNA that is present in a specific sample

## Principle

- ❑ The abundance of specific mRNA transcripts in a biological sample is a reflection of the expression levels of the corresponding genes (Manning et al., 2007). Information obtained from microarrays

## Application

- ❑ To associate differences in mRNA mixtures originating from different groups of individuals to phenotypic differences between the groups (Nachatomy et al., 2007)

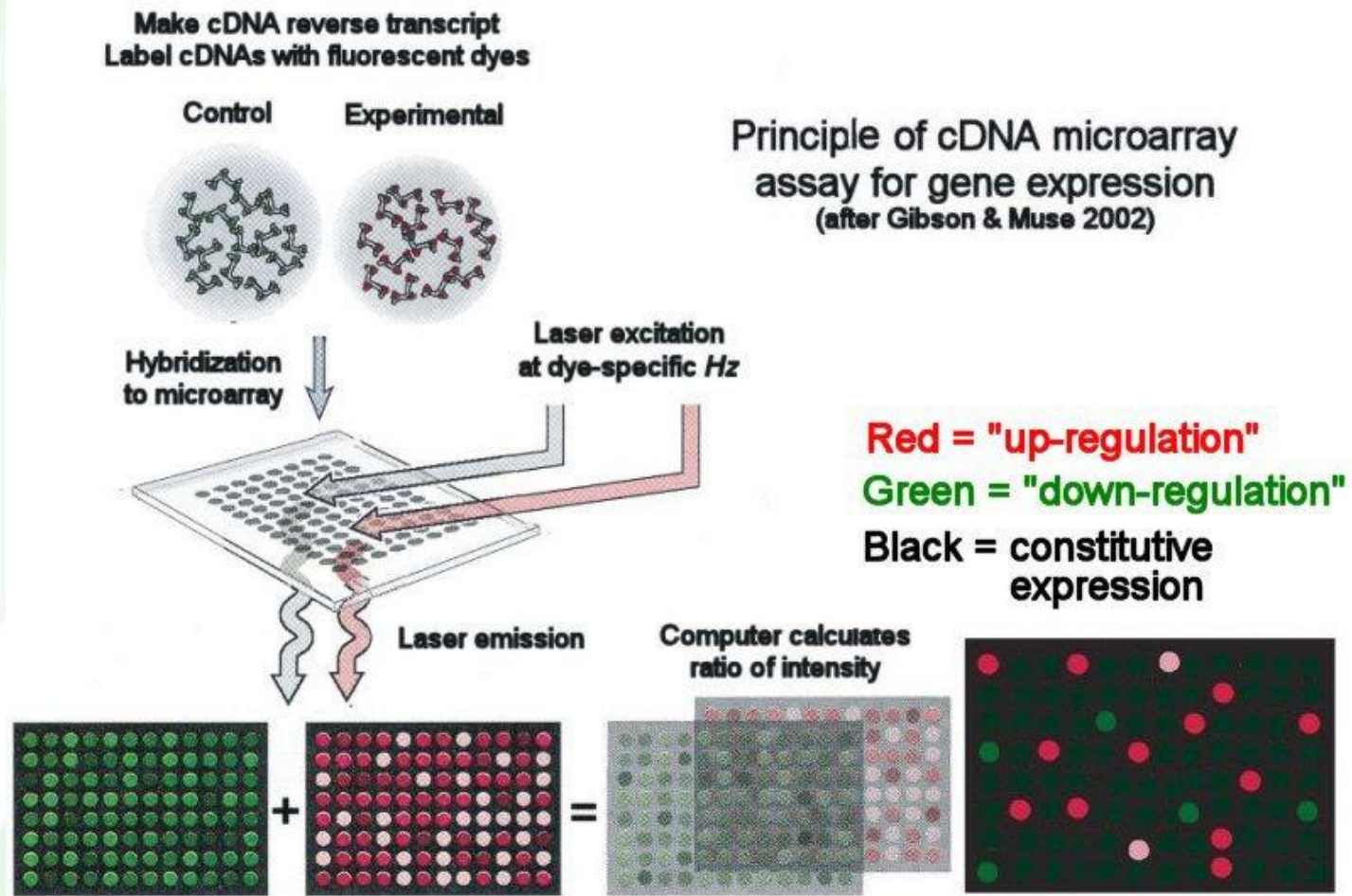
## Challenge

- ❑ The transcriptome in contrast to the genome is highly variable over time, between cell types and environmental changes (Celis et al., 2000).

## Recent improvement transcriptomics

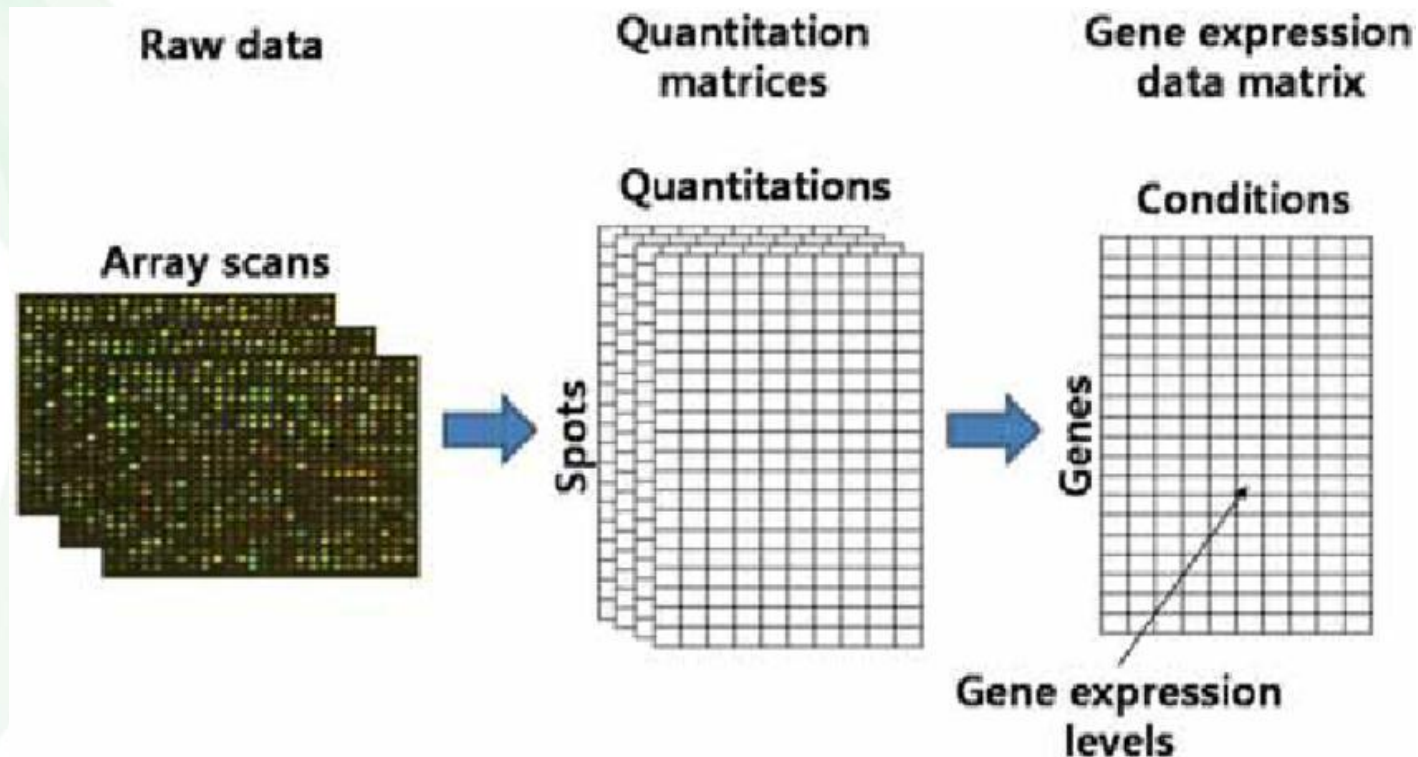
- ❑ RNAseq

# Transcriptomics





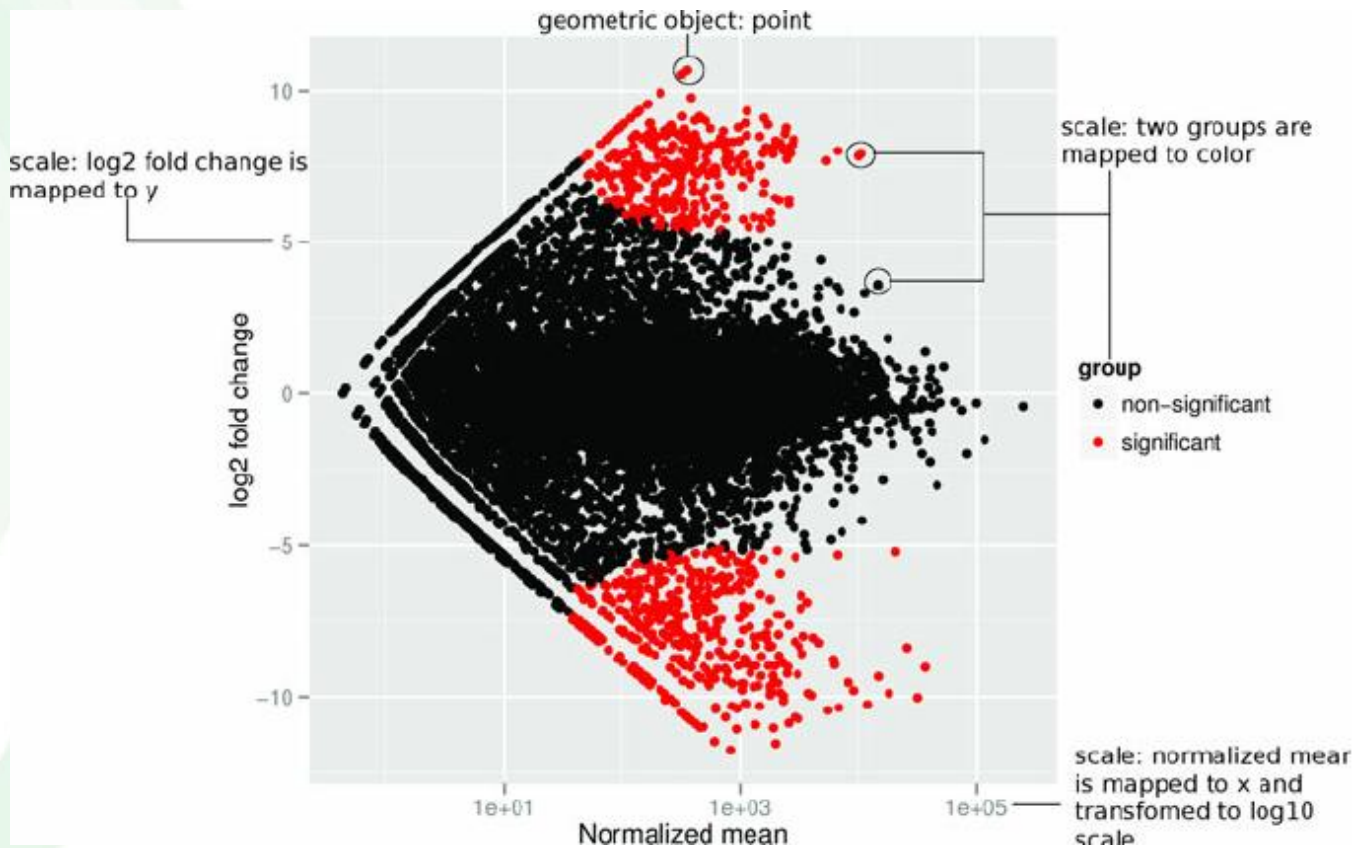
# Transcriptomics



**Statistical analysis:** linear models (one per gene) + empirical Bayes (**limma**)  
**Multiple comparisons:** FDR

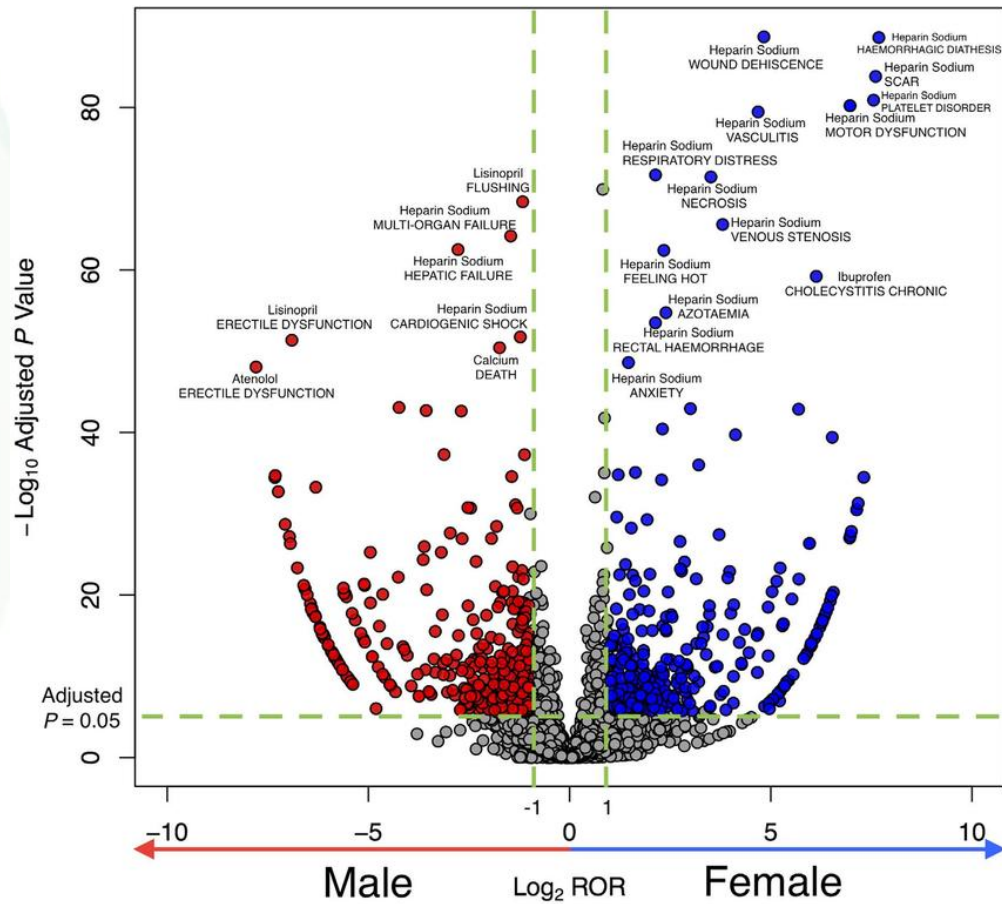
# Transcriptomics

## MA plot



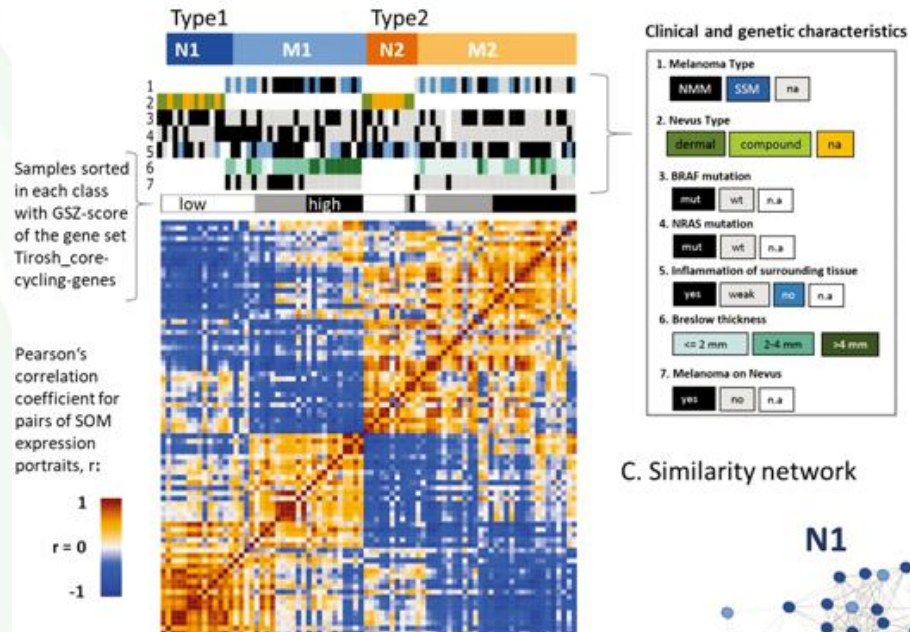
# Transcriptomics

## Volcano plot

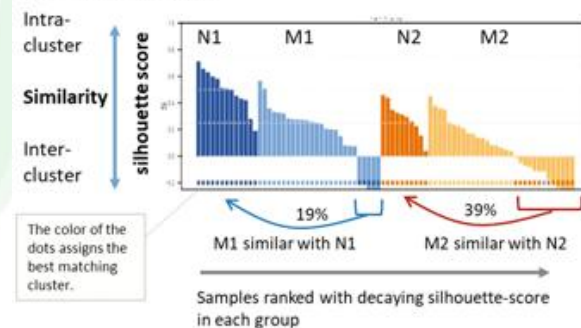


# Transcriptomics

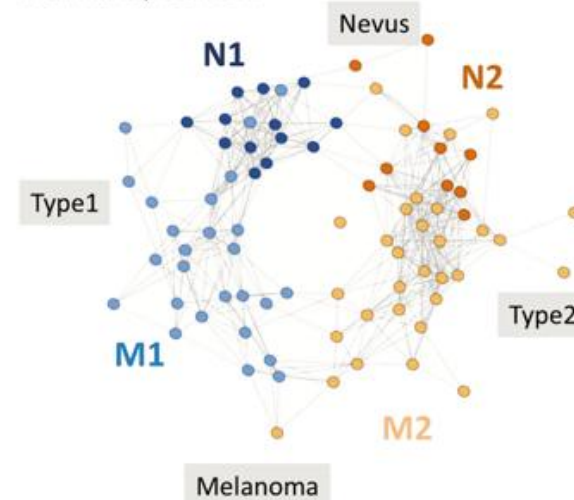
A. Pairwise correlation heatmap



B. Silhouette plot



C. Similarity network



# Epigenomics

## Epigenetic processes

- ☐ Mechanisms other than changes in DNA sequence that cause effect in gene transcription and gene silencing.
- ☐ Number of mechanisms of epigenomics but is mainly based on two mechanisms: DNA methylation and histone modification

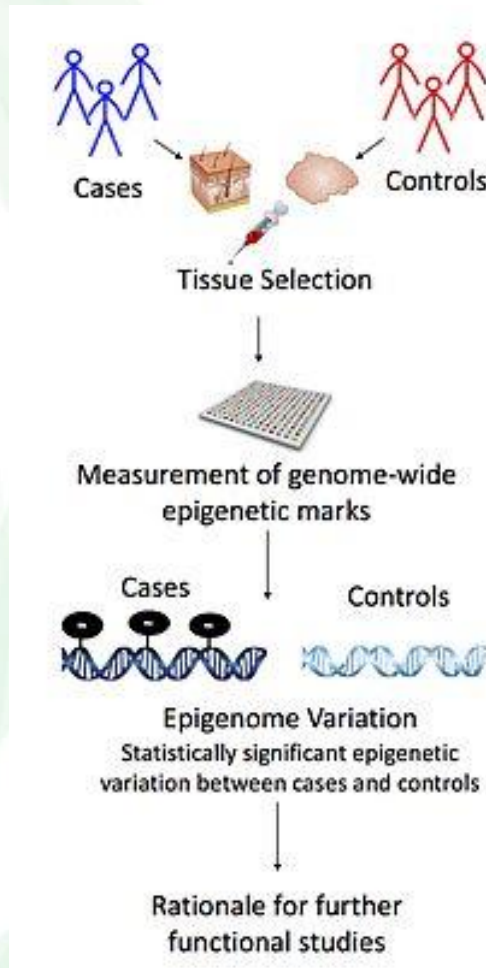
## Goal

- ☐ To study epigenetic processes on a large (ultimately genome-wide) scale to assess the effect on **disease** (450K and EPIC Illumina arrays)

## Association with disease

- ☐ Hypermethylation of CpG islands located in promoter regions of genes is related to gene silencing
- ☐ Altered gene silencing plays a causal role in human disease
- ☐ Histone proteins are involved in the structural packaging of DNA in the chromatin complex. Post translational histone modifications such as acetylation and methylation are believed to regulate chromatin structure and therefore gene expression

# Epigenomics



- ❑ 450K / EPIC Illumina arrays (**CpGs**)
- ❑ Statistical analyses using **linear models** (**robust regression** or beta regression)
- ❑ Multiple comparisons: **FDR**



# Proteomics

## Epigenetic processes

- ☐ Role proteins in biological systems. It consists of all proteins present in specific cell types or tissue and highly variable over time, between cell types and will change in response to changes in its environment (Fliser et al., 2007).
- ☐ The overall function of cells can be described by the proteins (intra- and intercellular and the abundance of these proteins (Sellers et al., 2003)
- ☐ Although all proteins are directly correlated to mRNA (transcriptome) , post translational modifications (PTM) and environmental interactions impede to predict from gene expression analysis alone (Hanash et al., 2008)

## Tools for proteomics

- ☐ Mass spectrometry (MS)
- ☐ Protein microarrays using capturing agents such as antibodies.

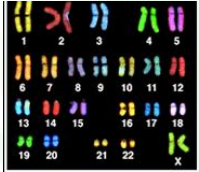
## Major focuses

- ☐ the identification of proteins and proteins interacting in protein-complexes
- ☐ Then the quantification of the protein abundance. The abundance of a specific protein is related to its role in cell function (Fliser et al., 2007)

# Metabolomics

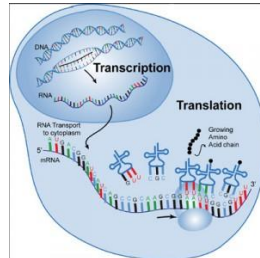
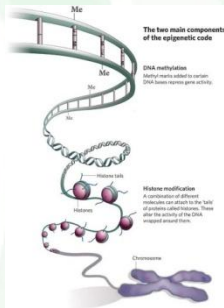
- ☐ The metabolome consists of small molecules (e.g. lipids or vitamins) that are also known as metabolites (Claudino et al., 2007).
- ☐ Metabolites are involved in the energy transmission in cells (metabolism) by interacting with other biological molecules following metabolic pathways.
- ☐ Metabolic phenotypes are the by-products of interactions between genetic, environmental, lifestyle and other factors (Holmes et al., 2008).
- ☐ The metabolome is highly variable and time dependent, and it consists of a wide range of chemical structures.
- ☐ An important challenge of metabolomics is to acquire qualitative and quantitative information with perturbation of environment (Jelly et al., 2010)
- ☐ Recent research on linking the metabolome with complex diseases or changes in environmental exposures

# Omics



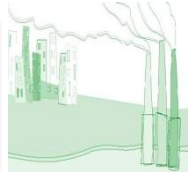
## GENOME

hereditary information (DNA)  
stable  
>99% equal between individuals  
1.5% coding genes



## EXPOSOME

dynamic  
diet, metals, air pollution, stress...

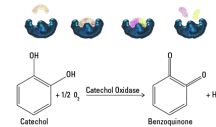


## EPIGENOME

changes in gene expression caused by mechanisms other than DNA sequence  
tissue and time specific

## TRANSCRIPTOME

gene expression (RNA)  
tissue and time specific



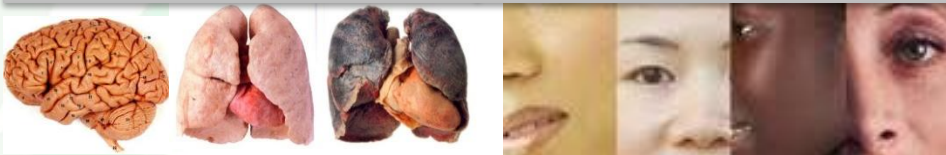
## PROTEOME

tissue and time specific

## METABOLOME

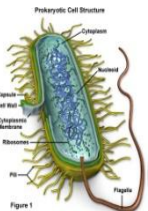
tissue and time specific

## DISEASOME (PHENOTYPE)



## METAGENOME

(metatranscriptome, virome...)  
bacteria and virus  
1-3% body's mass  
trillions of microorganisms



## Others

- ☐ Metagenome
- ☐ Immunome
- ☐ Pathogenome
- ☐ Regenome
- ☐ Psychogenome
- ☐ Phenome
- ☐ Exposome
- ☐ ...

# Omic data repositories

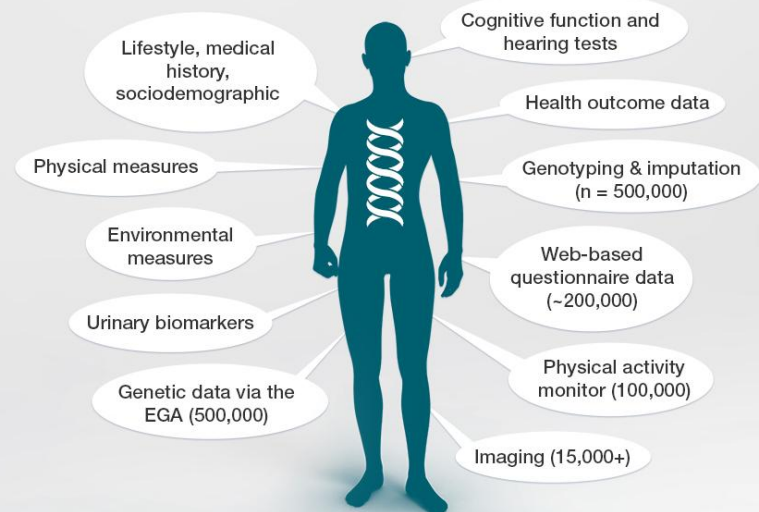
dbGaP (genomic data <https://www.ncbi.nlm.nih.gov/gap>)

EGA (genomic data <http://ega.crg.eu/>)



UK Biobank (genomic data <https://www.bdi.ox.ac.uk/>)

## Data on UK Biobank participants



# Omic data repositories

## GEO (transcriptomics – BioC: GEOquery)

<https://www.ncbi.nlm.nih.gov/geo/>

### Gene Expression Omnibus

GEO is a public functional genomics data repository supporting MIAME-compliant data submissions. Array- and sequence-based data are accepted. Tools are provided to help users query and download experiments and curated gene expression profiles.

Keyword or GEO Accession

#### Getting Started

[Overview](#)

[FAQ](#)

[About GEO DataSets](#)

[About GEO Profiles](#)

[About GEO2R Analysis](#)

[How to Construct a Query](#)

[How to Download Data](#)

#### Tools

[Search for Studies at GEO DataSets](#)

[Search for Gene Expression at GEO Profiles](#)

[Search GEO Documentation](#)

[Analyze a Study with GEO2R](#)

[Studies with Genome Data Viewer Tracks](#)

[Programmatic Access](#)

[FTP Site](#)

#### Browse Content

[Repository Browser](#)

DataSets: 4348

Series:  110236

Platforms: 19461

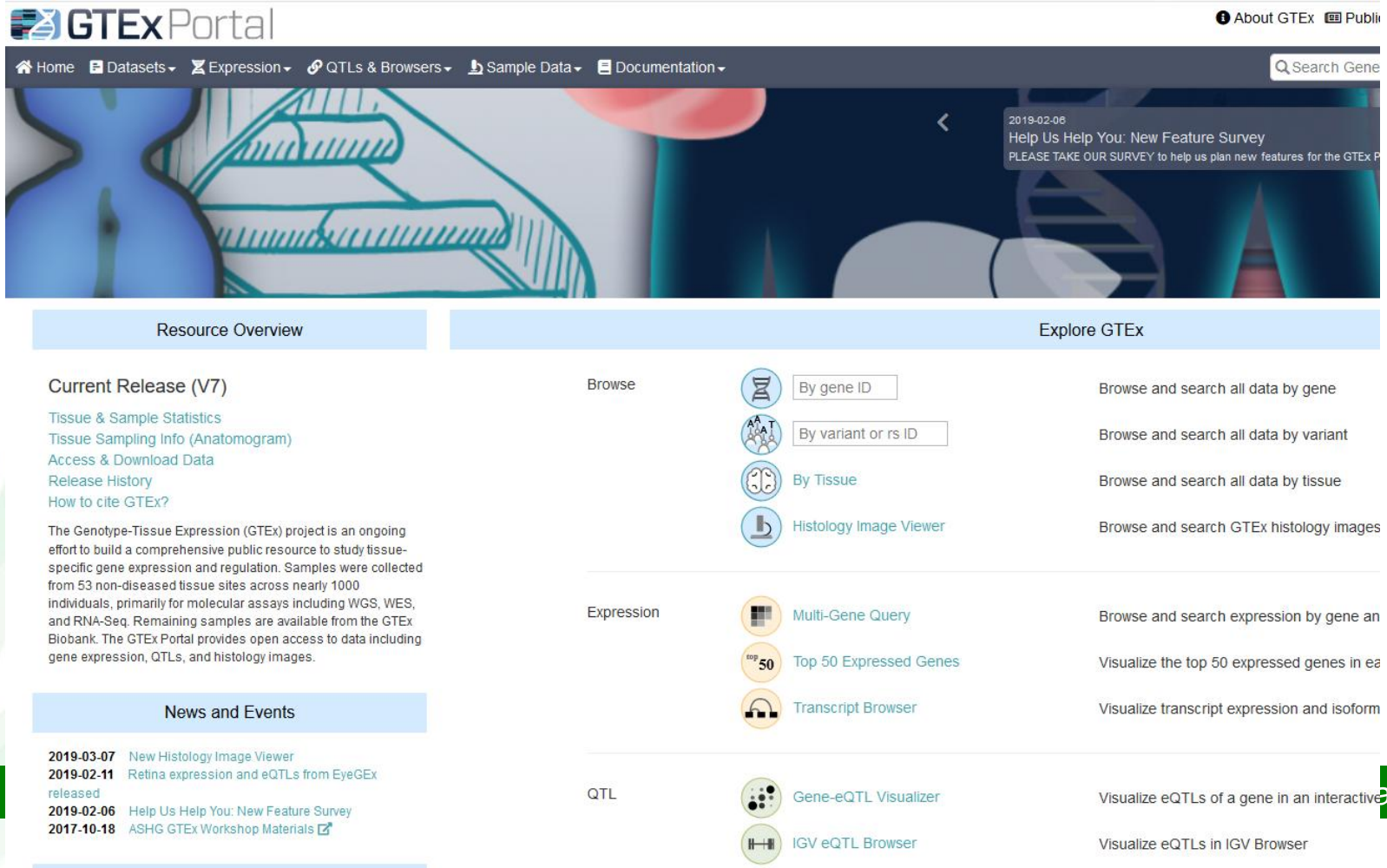
Samples: 2926596



# Omic data repositories

## GTeX(transcriptomics – genomics)

<https://gtexportal.org/home/>



The screenshot shows the GTEx Portal homepage. At the top is a green header with the text "Omic data repositories". Below this is the GTEx Portal logo and a navigation bar with links: Home, Datasets, Expression, QTLs & Browsers, Sample Data, and Documentation. A search bar is on the right. A banner image shows a stylized human figure with a DNA helix. A survey notification is visible: "2019-02-06 Help Us Help You: New Feature Survey PLEASE TAKE OUR SURVEY to help us plan new features for the GTEx P". The main content area is divided into two columns. The left column has a "Resource Overview" section with links for "Current Release (V7)", "Tissue & Sample Statistics", "Tissue Sampling Info (Anatomogram)", "Access & Download Data", "Release History", and "How to cite GTEx?". Below this is a "News and Events" section with dates and links for "New Histology Image Viewer", "Retina expression and eQTLs from EyeGE", "Help Us Help You: New Feature Survey", and "ASHG GTEx Workshop Materials". The right column has an "Explore GTEx" section with four categories: "Browse" (By gene ID, By variant or rs ID, By Tissue, Histology Image Viewer), "Expression" (Multi-Gene Query, Top 50 Expressed Genes, Transcript Browser), and "QTL" (Gene-eQTL Visualizer, IGV eQTL Browser). Each category has a brief description of the tool.

**Resource Overview**

**Current Release (V7)**

- [Tissue & Sample Statistics](#)
- [Tissue Sampling Info \(Anatomogram\)](#)
- [Access & Download Data](#)
- [Release History](#)
- [How to cite GTEx?](#)

The Genotype-Tissue Expression (GTEx) project is an ongoing effort to build a comprehensive public resource to study tissue-specific gene expression and regulation. Samples were collected from 53 non-diseased tissue sites across nearly 1000 individuals, primarily for molecular assays including WGS, WES, and RNA-Seq. Remaining samples are available from the GTEx Biobank. The GTEx Portal provides open access to data including gene expression, QTLs, and histology images.

**News and Events**

- 2019-03-07** [New Histology Image Viewer](#)
- 2019-02-11** [Retina expression and eQTLs from EyeGE released](#)
- 2019-02-06** [Help Us Help You: New Feature Survey](#)
- 2017-10-18** [ASHG GTEx Workshop Materials](#)

**Explore GTEx**

**Browse**

- [By gene ID](#): Browse and search all data by gene
- [By variant or rs ID](#): Browse and search all data by variant
- [By Tissue](#): Browse and search all data by tissue
- [Histology Image Viewer](#): Browse and search GTEx histology images

**Expression**

- [Multi-Gene Query](#): Browse and search expression by gene and tissue
- [Top 50 Expressed Genes](#): Visualize the top 50 expressed genes in each tissue
- [Transcript Browser](#): Visualize transcript expression and isoform

**QTL**

- [Gene-eQTL Visualizer](#): Visualize eQTLs of a gene in an interactive plot
- [IGV eQTL Browser](#): Visualize eQTLs in IGV Browser

# Omic data repositories

## TCGA (multiomic – BioC: several)

<https://www.cancer.gov/>

## Gene Expression Omnibus

GEO is a public functional genomics data repository supporting MIAME-compliant data submissions. Array- and sequence-based data are accepted. Tools are provided to help users query and download experiments and curated gene expression profiles.

Keyword or GEO Accession

### Getting Started

[Overview](#)

[FAQ](#)

[About GEO DataSets](#)

[About GEO Profiles](#)

[About GEO2R Analysis](#)

[How to Construct a Query](#)

[How to Download Data](#)

### Tools

[Search for Studies at GEO DataSets](#)

[Search for Gene Expression at GEO Profiles](#)

[Search GEO Documentation](#)

[Analyze a Study with GEO2R](#)

[Studies with Genome Data Viewer Tracks](#)

[Programmatic Access](#)

[FTP Site](#)

### Browse Content

[Repository Browser](#)

[DataSets:](#) 4348

[Series:](#)  110236

[Platforms:](#) 19461

[Samples:](#) 2926596

## Omic data repositories

**recount** (transcriptomic – BioC: several)

<https://jhubiostatistics.shinyapps.io/recount/>

recount2: analysis-ready RNA-seq gene and exon counts datasets

Datasets

Popular datasets

GTEx

TCGA

Documentation

Download data with R

Accessing recount2 via SciServer

Transcript counts are now available thanks to the work of Fu et al, bioRxiv, 2018. Exon counts are now from dis further information.



A multi-experiment resource of analysis-ready RNA-seq

**recount2** is an online resource consisting of RNA-seq gene and exon counts as well as coverage data were processed with [Rail-RNA](#) as described in the recount2 paper and at [Nellore et al, Genome Biology, 2016](#) which created the coverage extracted phenotype data, which we provide in their raw formats as well as in RangedSummarizedExperiment R objects (described in the [Summr](#) file, which can be used with the [derfinder](#) Bioconductor package to perform annotation-agnostic differential expression analysis at the expressed RangedSummarizeExperiment objects, phenotype tables, sample bigWigs, mean bigWigs, and file information tables are ready to use and free data for a specific study. By taking care of several preprocessing steps and combining many datasets into one easily-accessible website, we m

### Main publication

- **Collado-Torres L, Nellore A**, Kammers K, Ellis SE, Taub MA, Hansen KD, Jaffe AE, Langmead B, Leek JT. [Reproducible RNA-seq analy](#)

# Omic data repositories

[Datasets](#)
[Popular datasets](#)
[GTEx](#)
[TCGA](#)
[Documentation](#)
[Download data with R](#)
[Accessing recount2 via SciServer](#)
[Contribute your data](#)

This tab shows the information for the TCGA project. Due to its size, we also provide ranged summarized experiment objects (RSE) by tissue at the gene and exon levels

Show  entries

accession	number of samples	species	abstract	gene	exon
<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>
TCGA	11284	human	<p>The Cancer Genome Atlas (TCGA) is a collaboration between the National Cancer Institute (NCI) and the National Human Genome Research Institute (NHGRI) that has generated comprehensive, multi-dimensional maps of the key genomic changes in 33 types of cancer. The TCGA dataset, comprising more than two petabytes of genomic data, has been made publically available, and this genomic information helps the cancer research community to improve the prevention, diagnosis, and treatment of cancer.</p>	<p>RSE v2 counts v2 RSE v1 counts v1 RSE by tissue (version 2): adrenal gland bile duct bladder bone marrow brain breast cervix colorectal esophagus eye head and neck kidney liver lung lymph nodes ovary pancreas pleura prostate skin soft tissue stomach testis thymus thyroid uterus RSE by tissue (version 1): adrenal gland bile duct bladder bone marrow brain breast cervix colorectal esophagus eye head and neck kidney liver lung lymph nodes ovary pancreas pleura prostate skin soft tissue stomach testis thymus thyroid uterus</p>	<p>RSE v2 counts v2 RSE v1 counts v1 RSE by tissue (version 2): adrenal gland bile duct bladder bone marrow brain breast cervix colorectal esophagus eye head and neck kidney liver lung lymph nodes ovary pancreas pleura prostate skin soft tissue stomach testis thymus thyroid uterus RSE by tissue (version 1): adrenal gland bile duct bladder bone marrow brain breast cervix colorectal esophagus eye head and neck kidney liver lung lymph nodes ovary pancreas pleura prostate skin soft tissue stomach testis thymus thyroid uterus</p>

# Bioconductor

<https://www.bioconductor.org/>



Search:

**Home**

**Install**

**Help**

**Developers**

## About *Bioconductor*

Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data.

Bioconductor uses the R statistical programming language, and is open source and open development. It has two releases each year, and an active user community. Bioconductor is also available as an [AMI](#) (Amazon Machine Image) and a series of [Docker](#) images.

### Install »

- Discover [1649 software packages](#) available in *Bioconductor* release 3.8.

Get started with *Bioconductor*

- [Install Bioconductor](#)
- [Get support](#)
- [Latest newsletter](#)
- [Follow us on twitter](#)
- [Install R](#)

### Learn »

Master *Bioconductor* tools

- [Courses](#)
- [Support site](#)
- [Package vignettes](#)
- [Literature citations](#)
- [Common work flows](#)
- [FAQ](#)
- [Community resources](#)
- [Videos](#)