

What transcriptomic data tell us about

Nuno ~~disease~~ Barbosa Morais

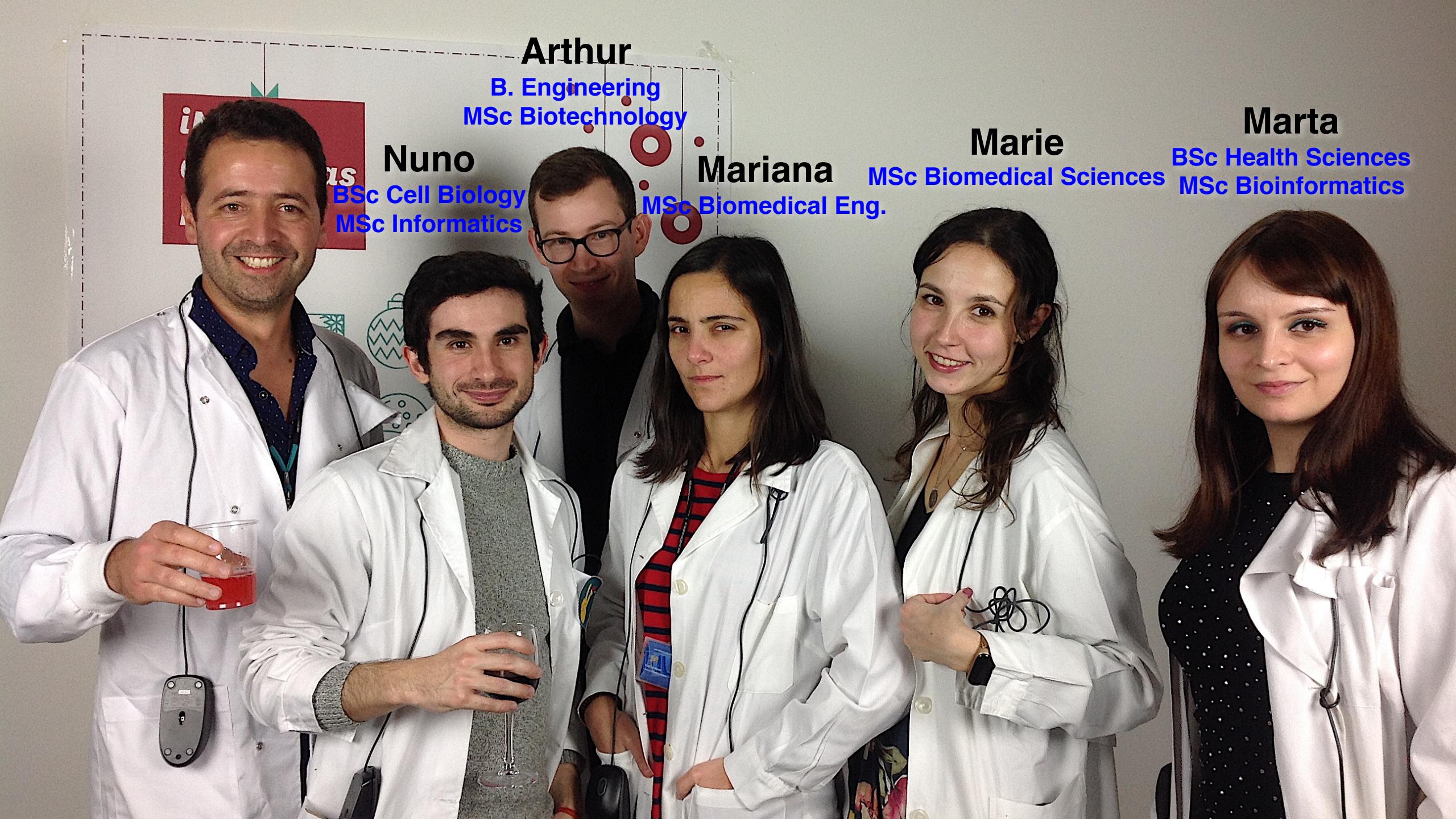




Instituto
de Medicina
Molecular

João
Lobo
Antunes





Arthur

B. Engineering
MSc Biotechnology

Nuno

BSc Cell Biology
MSc Informatics

O

Mariana

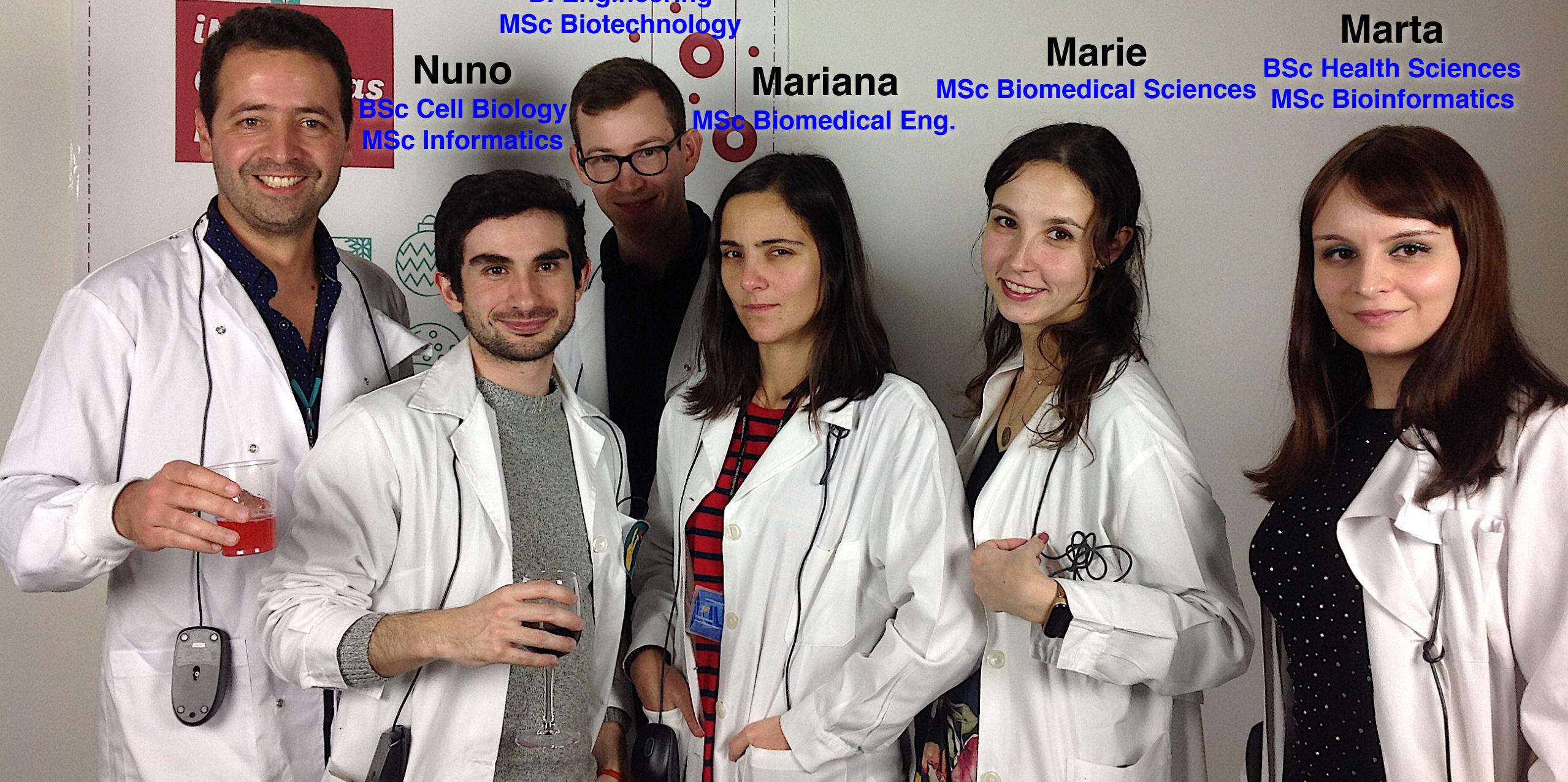
MSc Biomedical Eng.

Marta

BSc Health Sciences
MSc Bioinformatics

Marie

MSc Biomedical Sciences

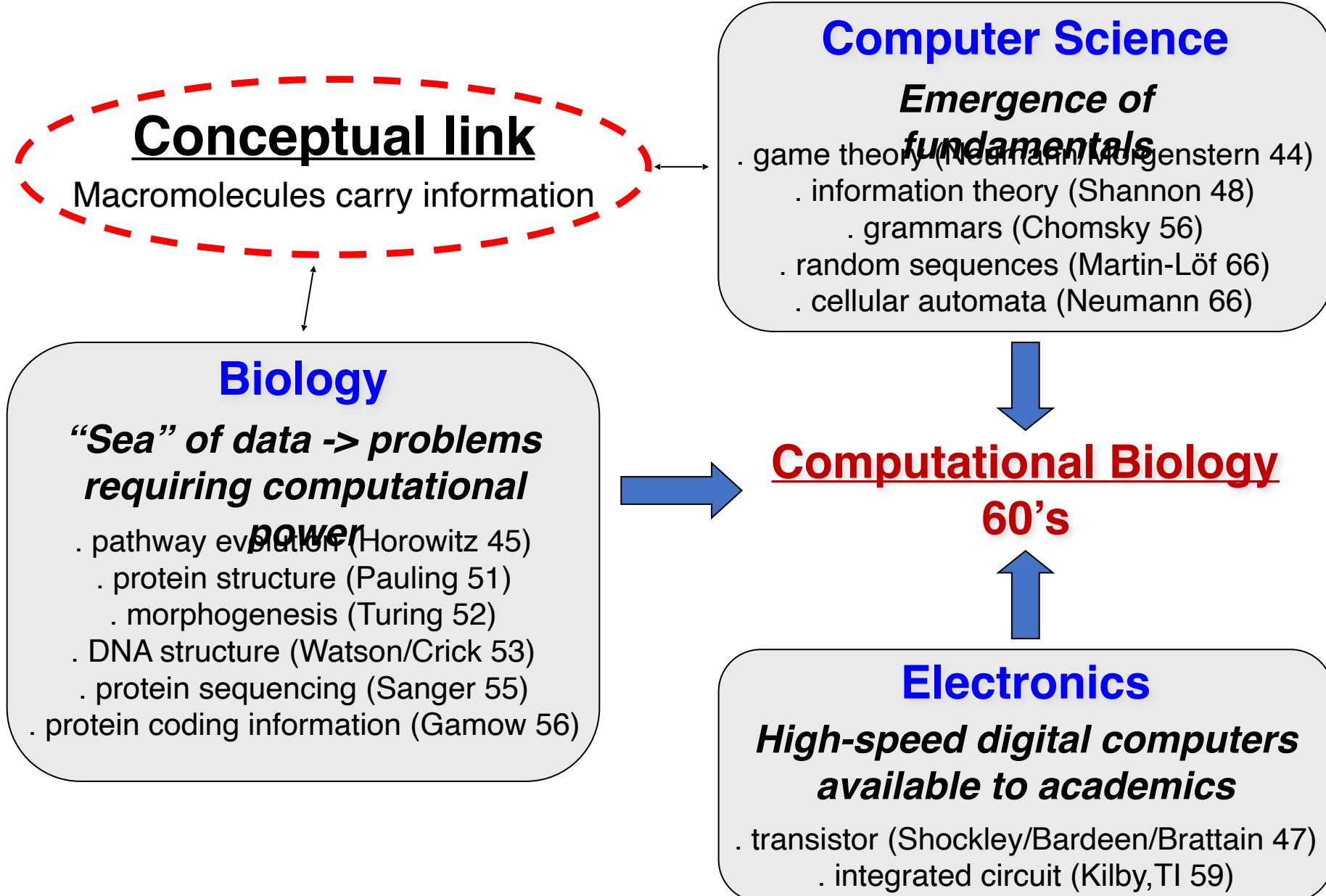


Outline

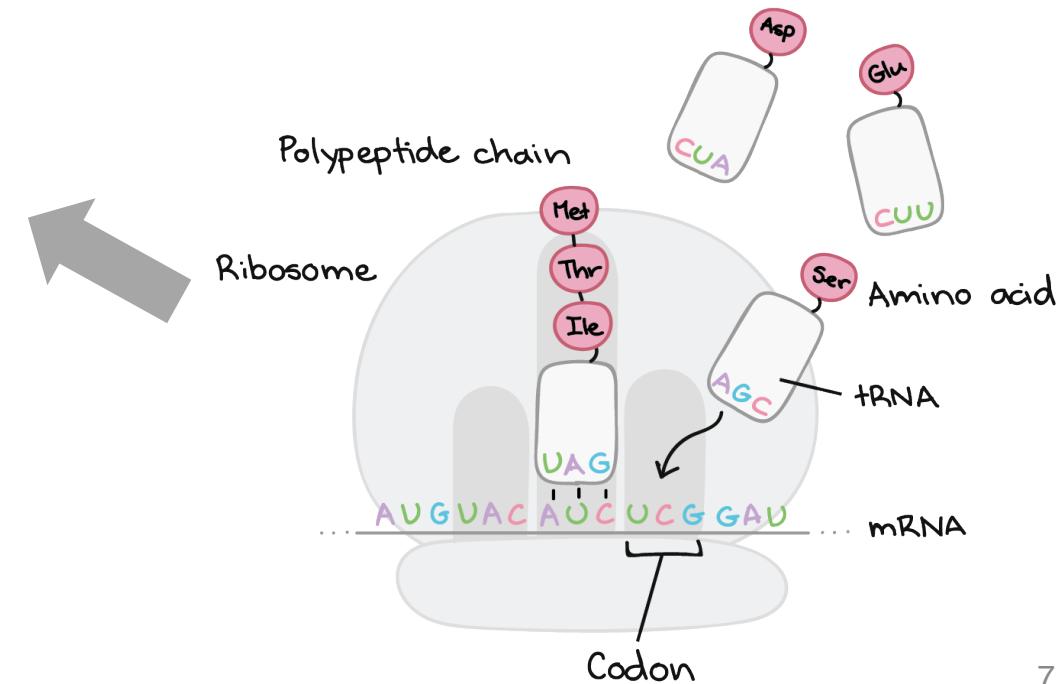
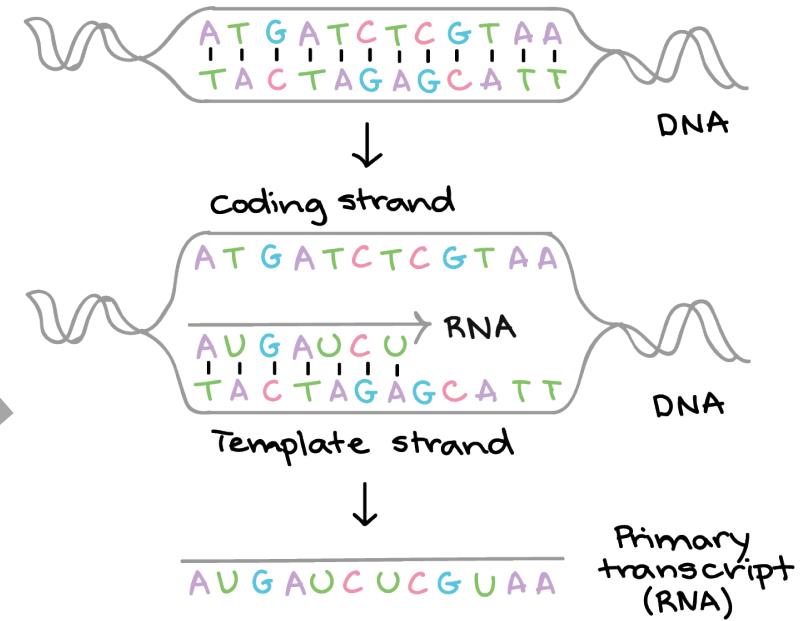
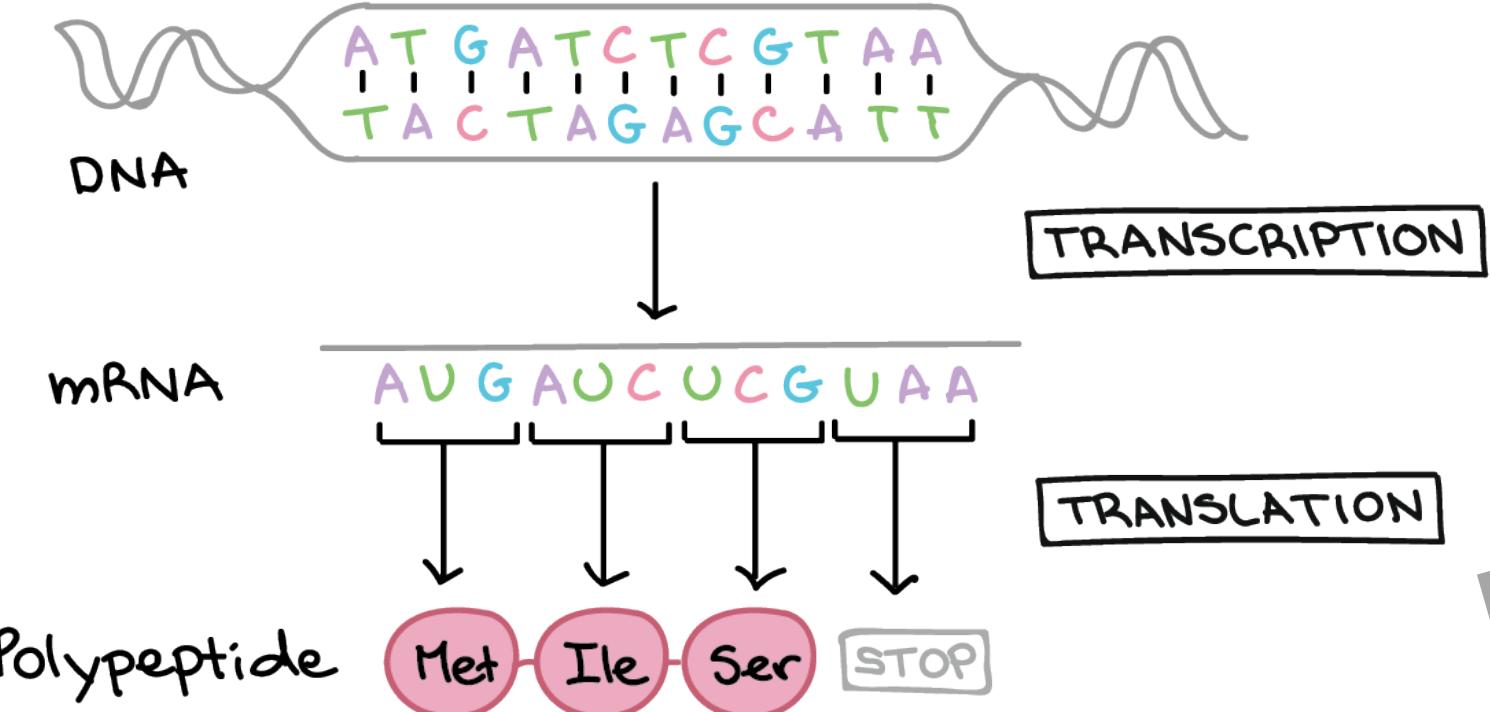
1. Biological macromolecules as carriers of biomedically relevant data
2. What we do: Disease Transcriptomics
3. Why biology needs more data scientists

1. Biological macromolecules as carriers of biomedically relevant data



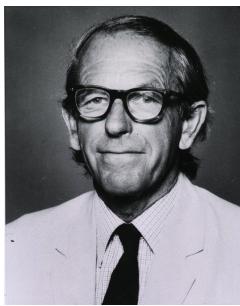


THE CENTRAL DOGMA



How to read DNA?

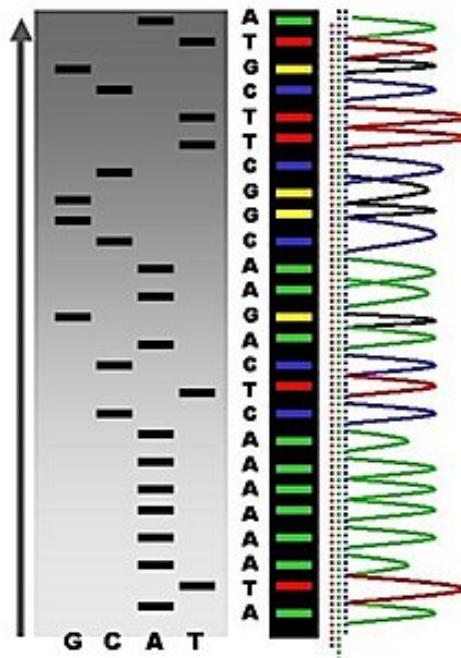
Fred Sanger
(1918-2013)



Chemistry 1958
Chemistry 1980

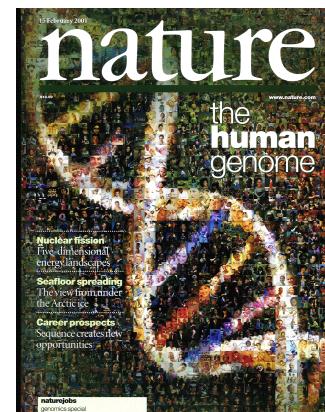
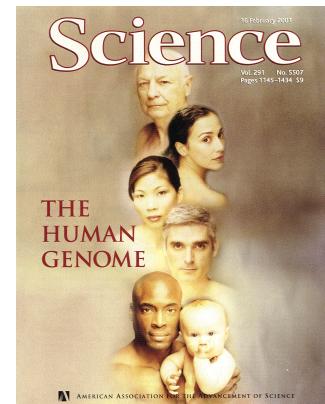
1977: Sanger sequencing

- *selective incorporation of chain-terminating dideoxynucleotides by DNA polymerase during in vitro DNA replication*
- state-of-the-art for >30 years



2001: first reference human genome

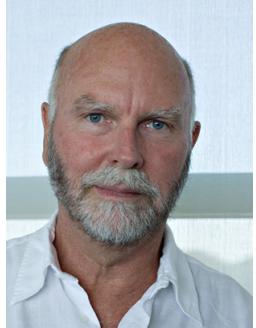
- ~\$1.000.000.000
- 10 years



Cost per Human Genome

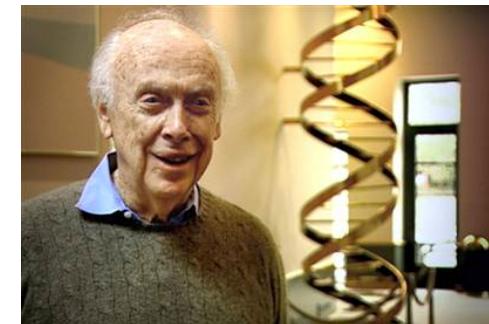


Individual genomics



**2007: Craig Venter's genome; ~10 M\$ and 1 year with the
new high-throughput sequencing technologies**

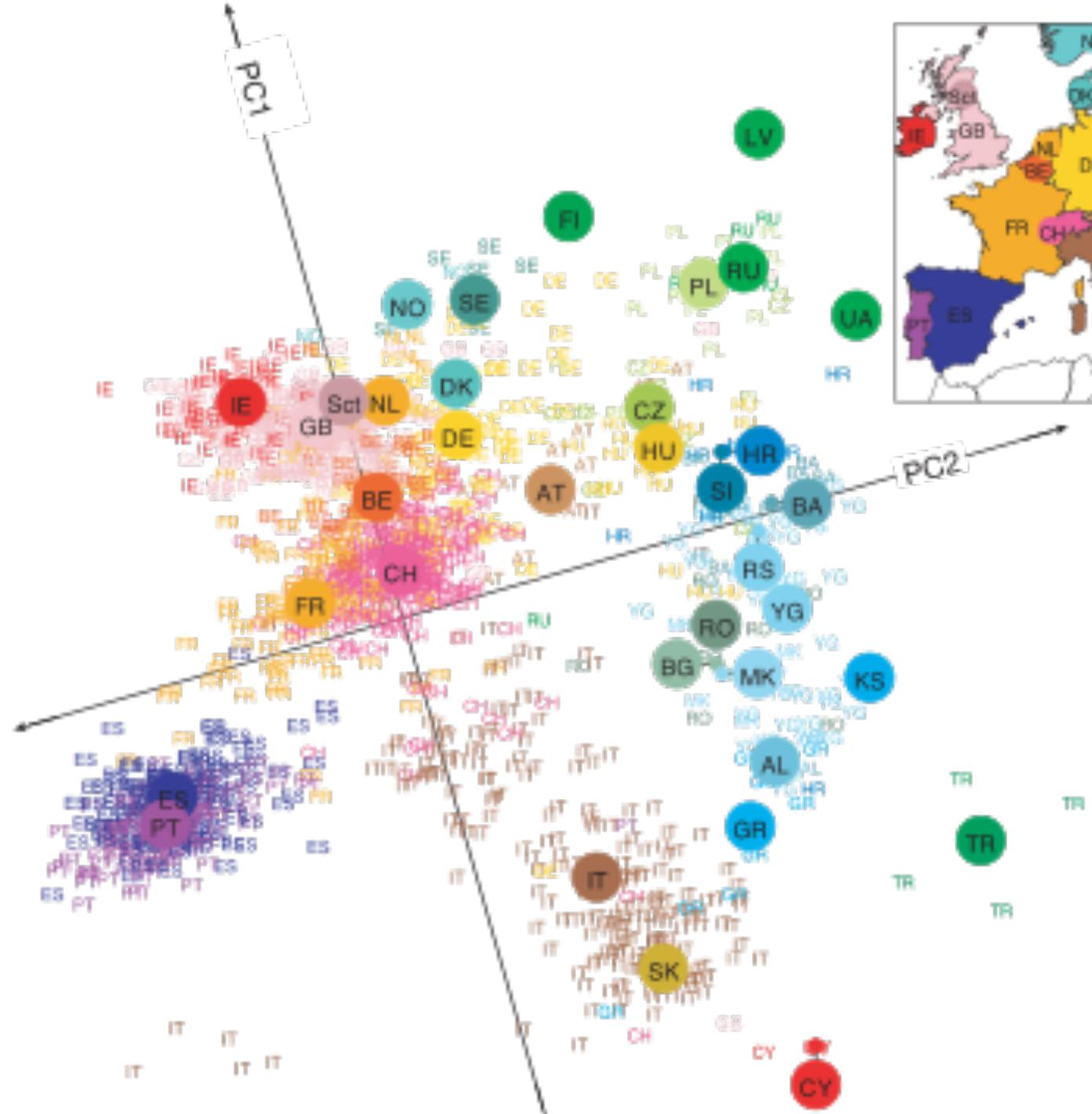
**2008: James Watson's genome;
~\$1.500.000 and 4 months**



2018: BROAD INSTITUTE sequences its 100Kth human genome, global total approaching 1M

Currently:
> 1Tb in <4 days





**Dimension reduction
and smart plots can
help us learn from
genomic data!**

nature

Vol 456 | 6 November 2008 | doi:10.1038/nature07331

LETTERS

Genes mirror geography within Europe

John Novembre^{1,2}, Toby Johnson^{4,5,6}, Katarzyna Bryc⁷, Zoltán Kutalik^{4,6}, Adam R. Boyko⁷, Adam Auton⁷, Amit Indap⁷, Karen S. King⁸, Sven Bergmann^{4,6}, Matthew R. Nelson⁸, Matthew Stephens^{2,3} & Carlos D. Bustamante⁷

Challenge: how to interpret the human genome?

- Knowing the sequence of ~ 3×10^9 bases (“alphabet” of only 4: ACGT) is not enough to understand how cells work
- All the somatic cells of an individual have the same genome and are quite diverse (different organs)



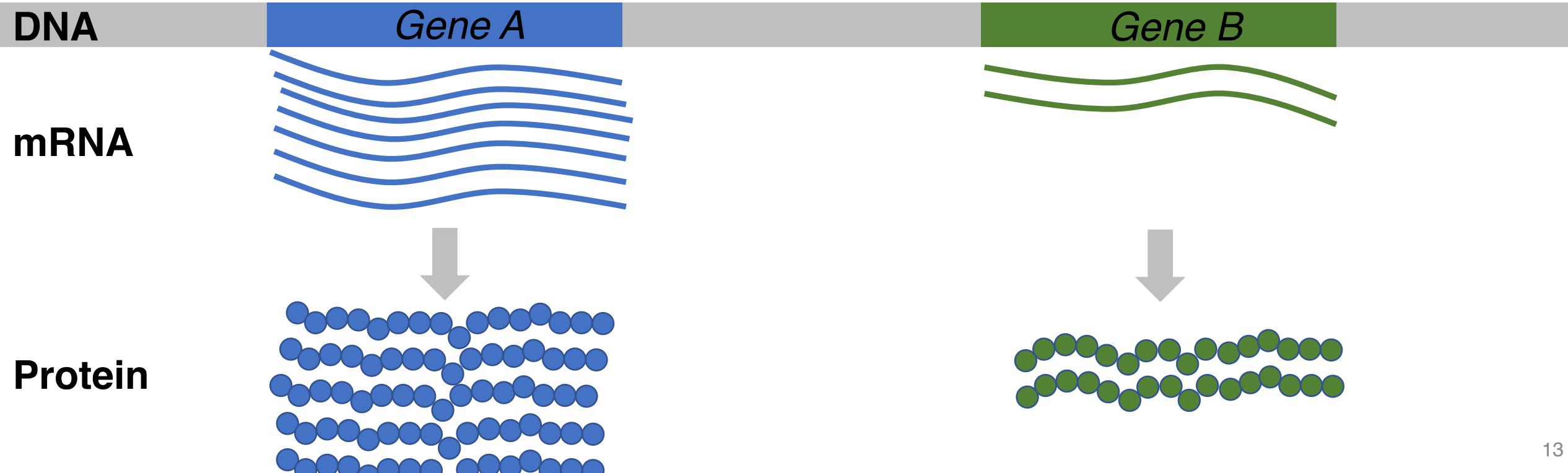
Post-genomic era: sequence → function

An approach: profiling the *expression* of all (1000s of) genes in parallel

Ultimate goal: measuring protein levels (possible but hard to do in a direct way)

Reasonable strategy: measuring levels of mRNA as surrogate for protein expression

↳ *sequence them (RNA-seq) and count them*



How the data look like (breast cancer)

Gene expression (counts)

Gene	TOGA-A1-A0SB-01	TOGA-A1-A0SD-01	TOGA-A1-A0SE-01	TOGA-A1-A0SF-01	TOGA-A1-A0SM-01
A1BG	164	546	1341	836	
A1CF	0	0	0	1	
A2BP1	22	1	2	0	
A2LD1	127	331	498	526	
A2ML1	94	144	114	77	
A2M	102123	107181	101192	50316	
A4GALT	890	1409	1711	1142	
A4GNT	6	5	2	2	
AAA1	0	0	1	0	
AAAS	2139	2219	4294	3052	
AACSL	2930	3	4	1	
AACS	6533	3102	5271	1248	
AADACL2	0	0	1	0	
AADACL3	0	0	0	0	
AADACL4	1	0	0	0	
AADAC	42	4	13	4	
AADAT	694	226	200	67	
AAGAB	2913	9069	10465	3366	
AAK1	3811	5545	6963	2341	
AAMP	6514	9517	18397	8569	
AANAT	9	2	11	0	
AARS2	1764	2193	3338	1463	
AARSD1	2354	1846	3391	2169	
AARS	10855	9400	16774	6233	
AASDHPPPT	3234	2936	3766	3134	
AASDH	1299	1589	2548	902	
AASS	1587	1852	2014	852	
AATF	2870	5656	9522	4625	
AATK	317	312	736	169	
ABAT	1422	3956	12660	3237	
ABC10	143	177	319	146	
ABCML1	271	142	761	120	

Sample annotation

Patient ID	Gender	Age	Menopausal	Ethnicity	Year of diagnosis	Anatomic subdivision	Surgical procedure
TCGA-A1-A0SB	female	70	post (prior bi)	not hispanic or latino	2008	left	lumpectomy
TCGA-A1-A0SD	female	59	NA	not hispanic or latino	2005	left	lumpectomy
TCGA-A1-A0SE	female	56	pre (<6 mont)	not hispanic or latino	2005	left upper outer quadrant	modified radical
TCGA-A1-A0SF	female	54	pre (<6 mont)	not hispanic or latino	2006	left	modified radical
TCGA-A1-A0SG	female	61	post (prior bi)	not hispanic or latino	2006	right	lumpectomy
TCGA-A1-A0SH	female	39	pre (<6 mont)	not hispanic or latino	2006	left upper inner quadrant	lumpectomy
TCGA-A1-A0SI	female	52	NA	not hispanic or latino	2007	right	lumpectomy
TCGA-A1-A0SJ	female	39	NA	not hispanic or latino	2006	left	modified radical
TCGA-A1-A0SK	female	54	indeterminate	not hispanic or latino	2007	right upper outer quadrant	lumpectomy
TCGA-A1-A0SM	male	77	NA	not hispanic or latino	2007	left	modified radical
TCGA-A1-A0SN	female	50	post (prior bi)	not hispanic or latino	2007	left	lumpectomy
TCGA-A1-A0SO	female	67	post (prior bi)	not hispanic or latino	2007	right	modified radical
TCGA-A1-A0SP	female	40	NA	not hispanic or latino	2007	right	lumpectomy
TCGA-A1-A0SQ	female	45	pre (<6 mont)	not hispanic or latino	2007	left	modified radical
TCGA-A2-A04N	female	66	post (prior bi)	hispanic or latino	2002	right	lumpectomy
TCGA-A2-A04P	female	36	pre (<6 mont)	NA	2003	left	lumpectomy
TCGA-A2-A04Q	female	48	post (prior bi)	NA	2004	right upper outer quadrant	simple mastectomy
TCGA-A2-A04R	female	36	pre (<6 mont)	not hispanic or latino	2004	left upper inner quadrant	lumpectomy
TCGA-A2-A04T	female	62	post (prior bi)	NA	2004	left upper inner quadrant	lumpectomy
TCGA-A2-A04U	female	47	peri (6-12 mont)	not hispanic or latino	2004	right	simple mastectomy
TCGA-A2-A04V	female	39	pre (<6 mont)	not hispanic or latino	2005	right lower outer quadrant	simple mastectomy
TCGA-A2-A04W	female	50	pre (<6 mont)	not hispanic or latino	2005	right upper outer quadrant	modified radical
TCGA-A2-A04X	female	34	pre (<6 mont)	NA	2006	left	simple mastectomy
TCGA-A2-A04Y	female	53	post (prior bi)	NA	2008	right upper outer quadrant	simple mastectomy
TCGA-A2-A0CL	female	37	pre (<6 mont)	not hispanic or latino	2006	right lower outer quadrant	modified radical
TCGA-A2-A0CM	female	40	pre (<6 mont)	not hispanic or latino	2003	left	simple mastectomy
TCGA-A2-A0CP	female	60	post (prior bi)	not hispanic or latino	2003	right	lumpectomy
TCGA-A2-A0CQ	female	62	post (prior bi)	not hispanic or latino	2003	right	lumpectomy
TCGA-A2-A0CS	female	73	post (prior bi)	not hispanic or latino	2004	left	modified radical
TCGA-A2-A0CT	female	71	post (prior bi)	not hispanic or latino	2005	right upper outer quadrant	simple mastectomy
TCGA-A2-A0CU	female	73	post (prior bi)	not hispanic or latino	2005	right lower outer quadrant	simple mastectomy
TCGA-A2-A0CV	female	41	post (prior bi)	not hispanic or latino	2005	left lower outer quadrant	modified radical

Biological samples

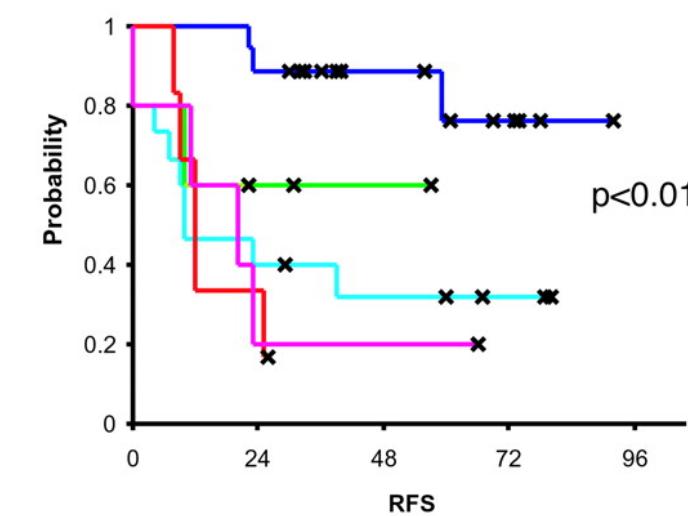
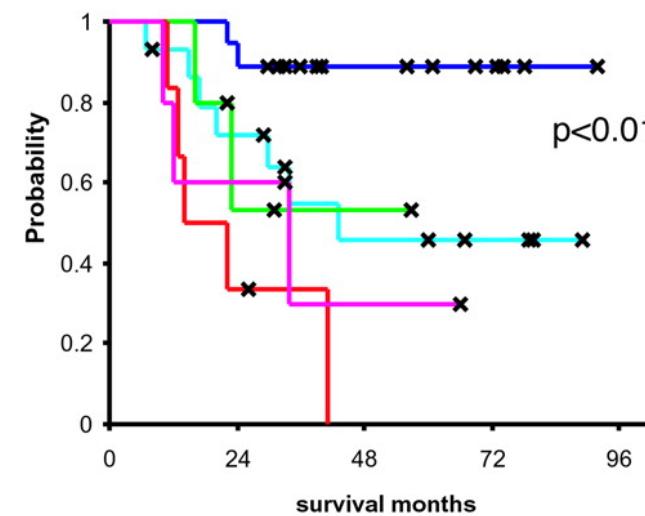
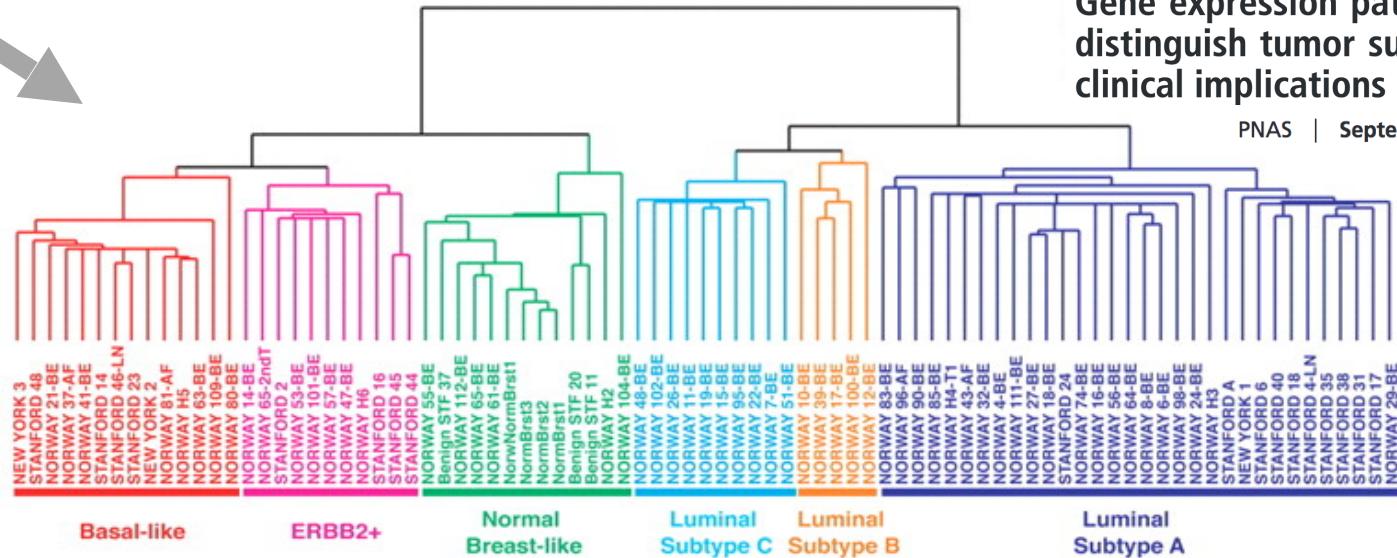
Samples



Breast cancer clustering and classification

Gene expression patterns of breast carcinomas
distinguish tumor subclasses with
clinical implications

PNAS | September 11, 2001 | vol. 98 | no. 19 | 10869–10874



X Censored, — Lum A, — Lum B+C, — NorB-like, — Basal, — ERBB2+

Breast cancer molecular subtyping in the clinic

Gene expression patterns of breast carcinomas
distinguish tumor subclasses with
clinical implications

Therese Sørlie^{a,b,c}, Charles M. Perou^{a,d}, Robert Tibshirani^e, Turid Aas^f, Stephanie Geisler^g, Hilde Johnsen^b, Trevor Hastie^e, Michael B. Eisen^h, Matt van de Rijnⁱ, Stefanie S. Jeffrey^j, Thor Thorsen^k, Hanne Quist^l, John C. Matese^c, Patrick O. Brown^m, David Botstein^c, Per Eystein Lønning^o, and Anne-Lise Børresen-Dale^{b,n}

PNAS | September 11, 2001 | vol. 98 | no. 19 | 10869–10874



Predictive Biomarkers and Personalized Medicine

Clinical
Cancer
Research

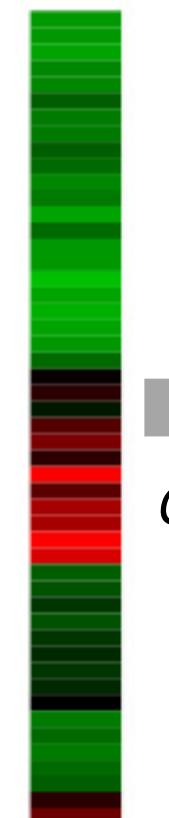
A 50-Gene Intrinsic Subtype Classifier for Prognosis and
Prediction of Benefit from Adjuvant Tamoxifen

Stephen K. Chia¹, Vivien H. Bramwell³, Dongsheng Tu⁴, Lois E. Shepherd⁴, Shan Jiang⁴, Tammi Vickery⁶,
Elaine Mardis⁶, Samuel Leung², Karen Ung², Kathleen I. Pritchard⁵, Joel S. Parker⁷, Philip S. Bernard⁸,
Charles M. Perou⁷, Matthew J. Ellis⁵, and Torsten O. Nielsen²



PAM50 (PROSIGNA)

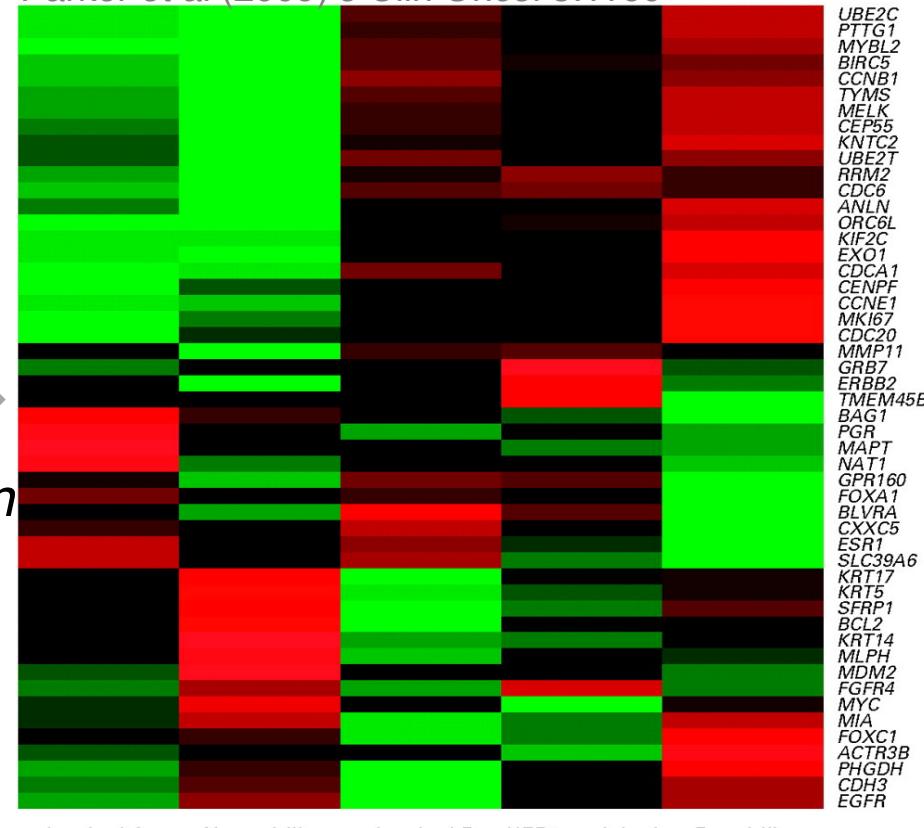
PAM50 (Prosigna®) is a tumor profiling test that helps determine the benefit of using chemotherapy in addition to hormone therapy for some estrogen receptor-positive (ER-positive), HER2-negative breast cancers.



Patient

Classification

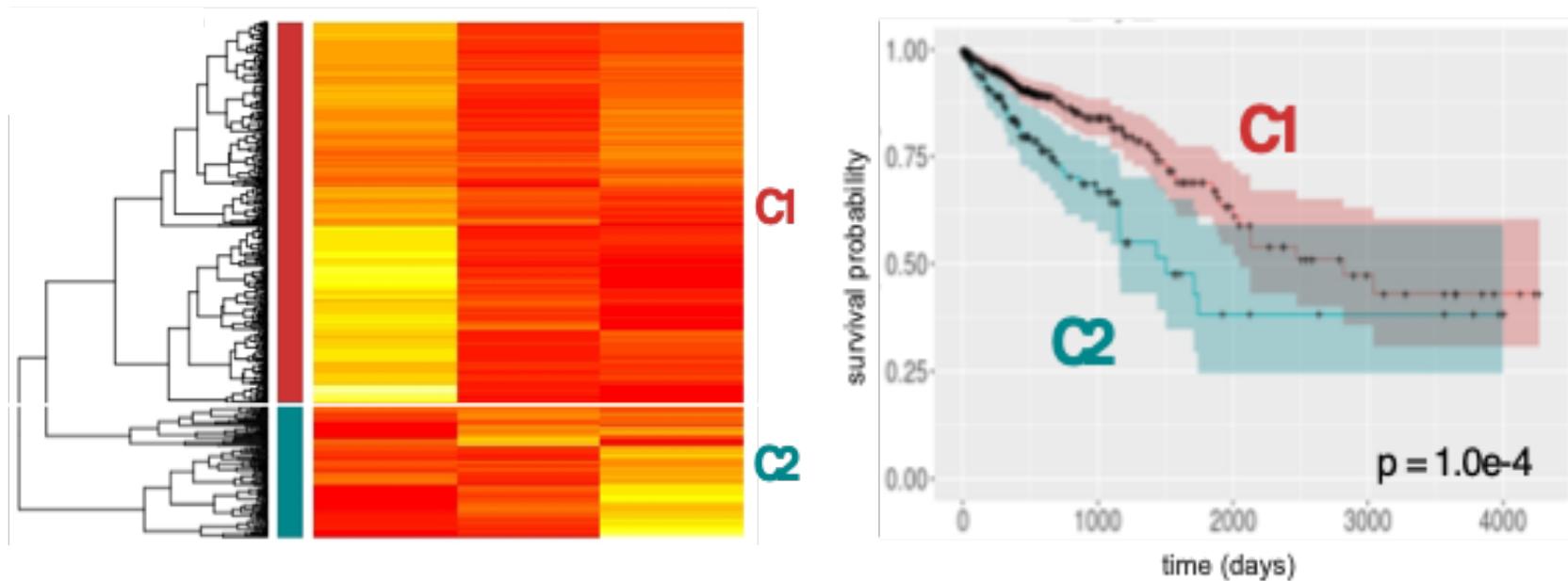
Parker et al (2009) J Clin Oncol 8:1160



Luminal A Normal-like Luminal B HER2-enriched Basal-like



2. What we do: *Disease Transcriptomics*



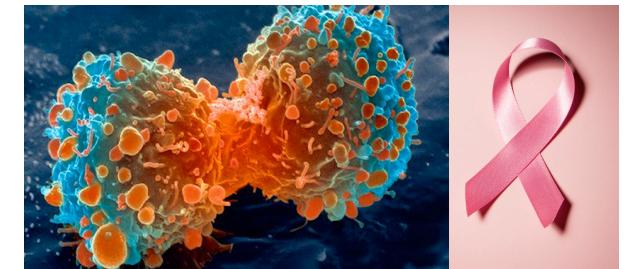
Goal

Understand how ageing-associated molecular (RNA) changes in human tissues increase proneness to disease

(Disease) models



*Neurodegeneratio
n*



Cancer

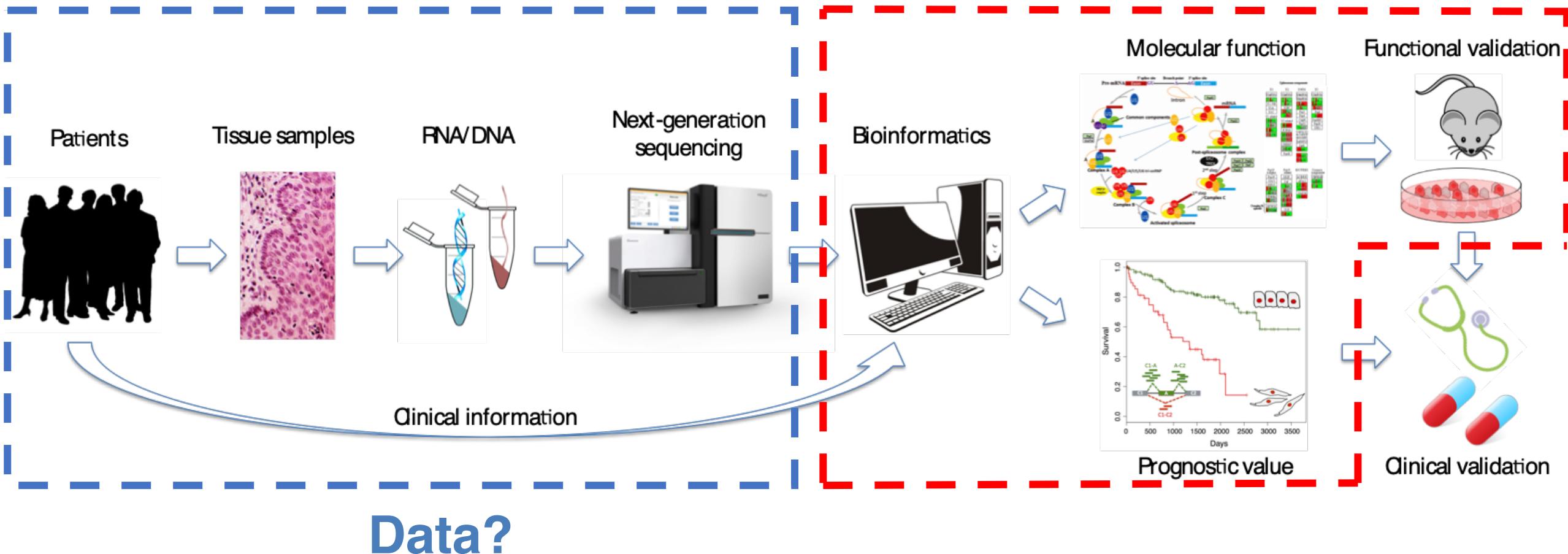
Ageing



Why the transcriptome?

(the set of all RNA molecules in one cell or a population of cells)

- Early and accurately profileable measure of cells' response to stimuli
- Our expertise: computational biology + transcriptional regulation
- Timely: recent availability of data



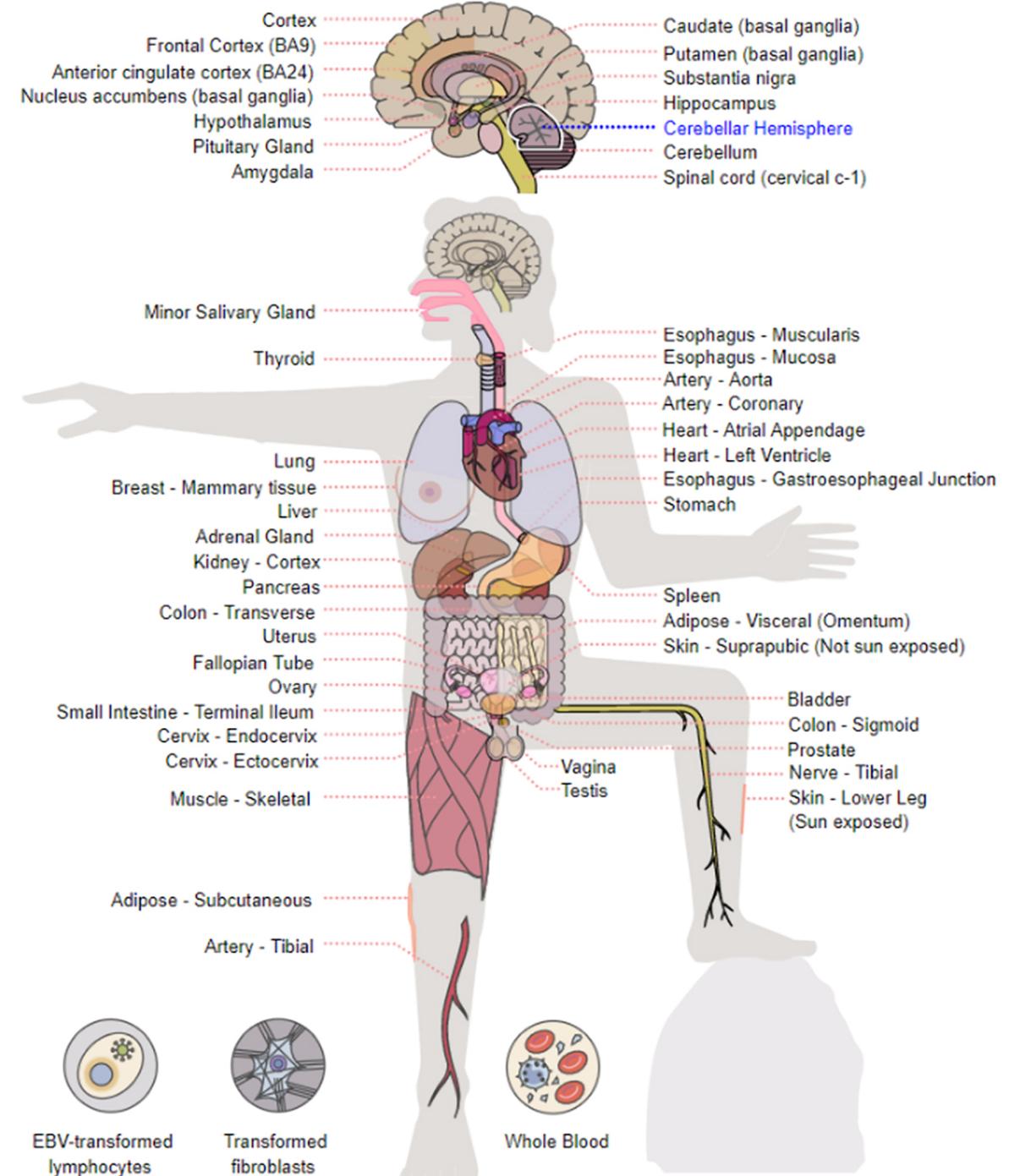
Big public data!



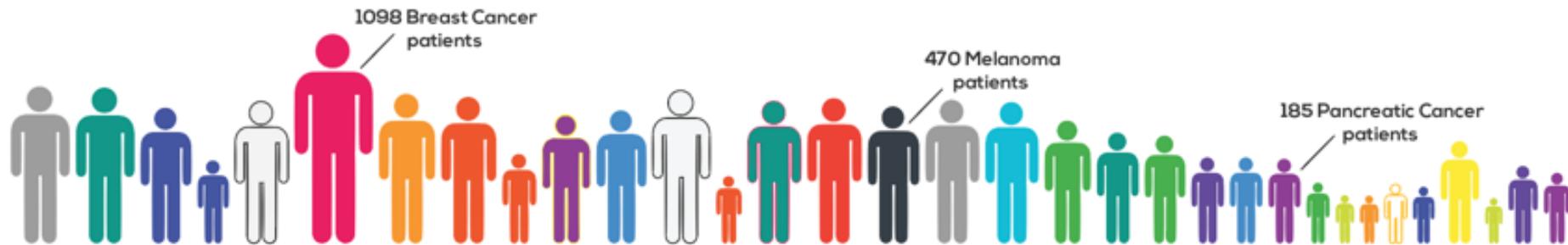
Genotype-Tissue Expression

RNA-seq from 53 human tissues

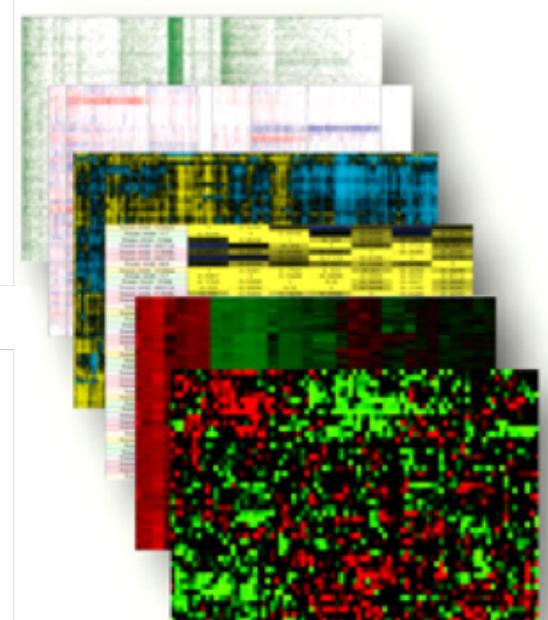
- >10.000 samples (100s per tissue)
- clinically annotated ([age](#), gender, ethnicity, cause of death, [diseases](#))



Matched tumor & normal tissues from more than **11,000** patients, representing **33** cancer types.



- █ Clinical ←
- █ SNP6 CopyNum
- █ LowPass DNaseq CopyNum
- █ Mutation Annotation File
- █ methylation
- █ miR
- █ miRSeq
- █ mRNA
- █ mRNASeq ←
- █ raw Mutation Annotation File
- █ Reverse Phase Protein Array





ConnectivityMap

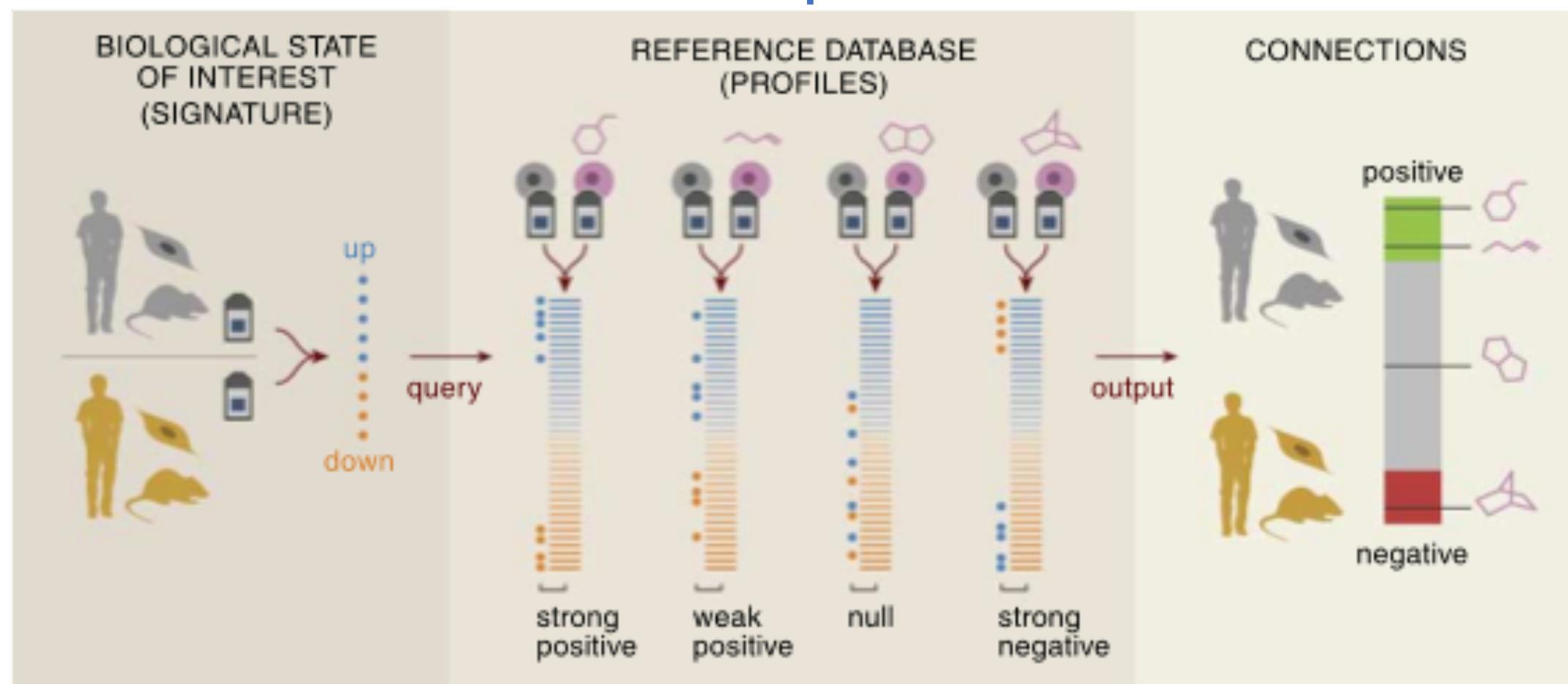
Unravel biology with the world's largest
perturbation-driven gene expression dataset.



Lamb *et al* (2006) *Science* 313:1929

Subramanian *et al* (2017) *Cell* 171:1437

Perturbations:
~20k small molecules
>5k genes

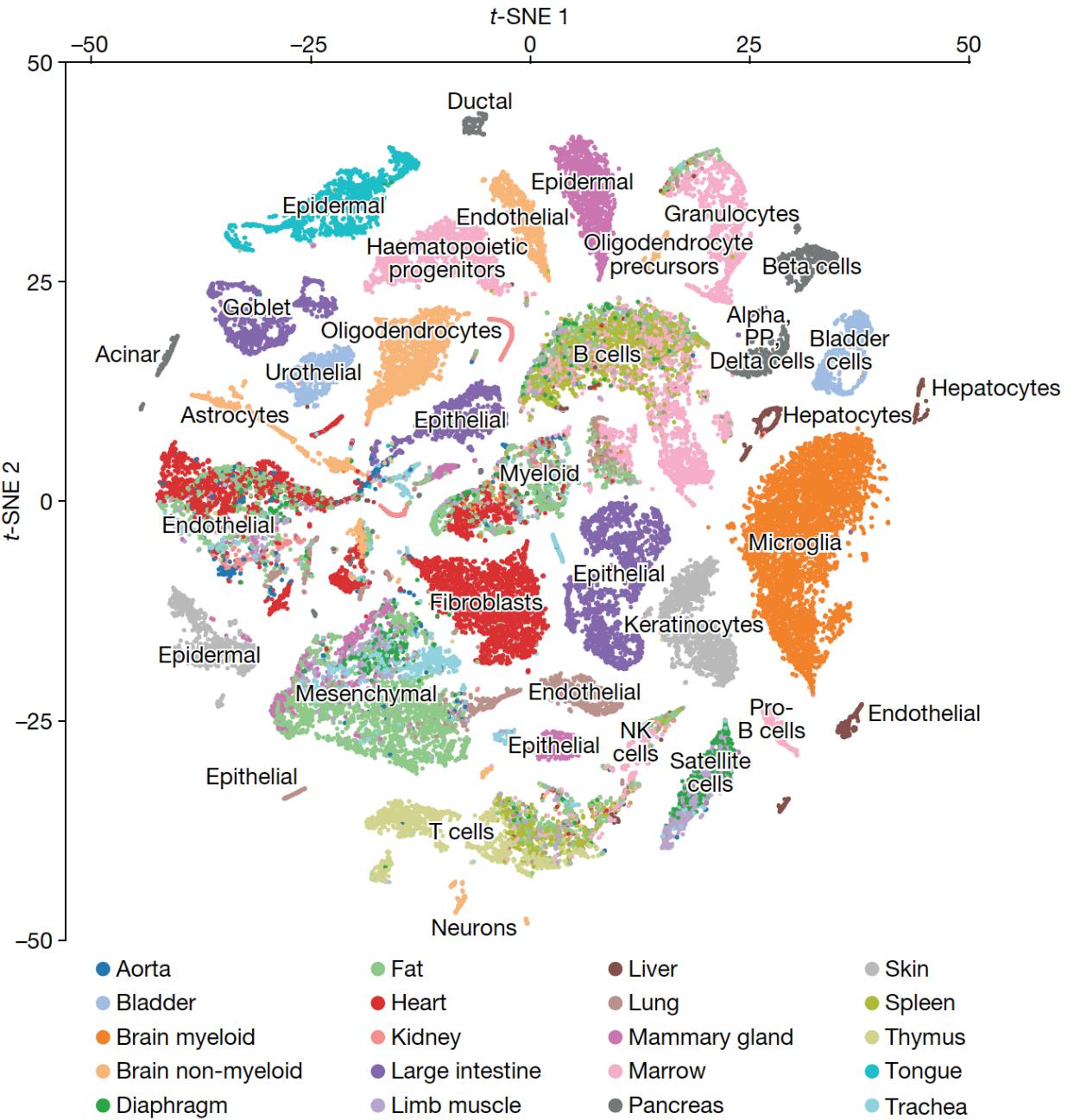
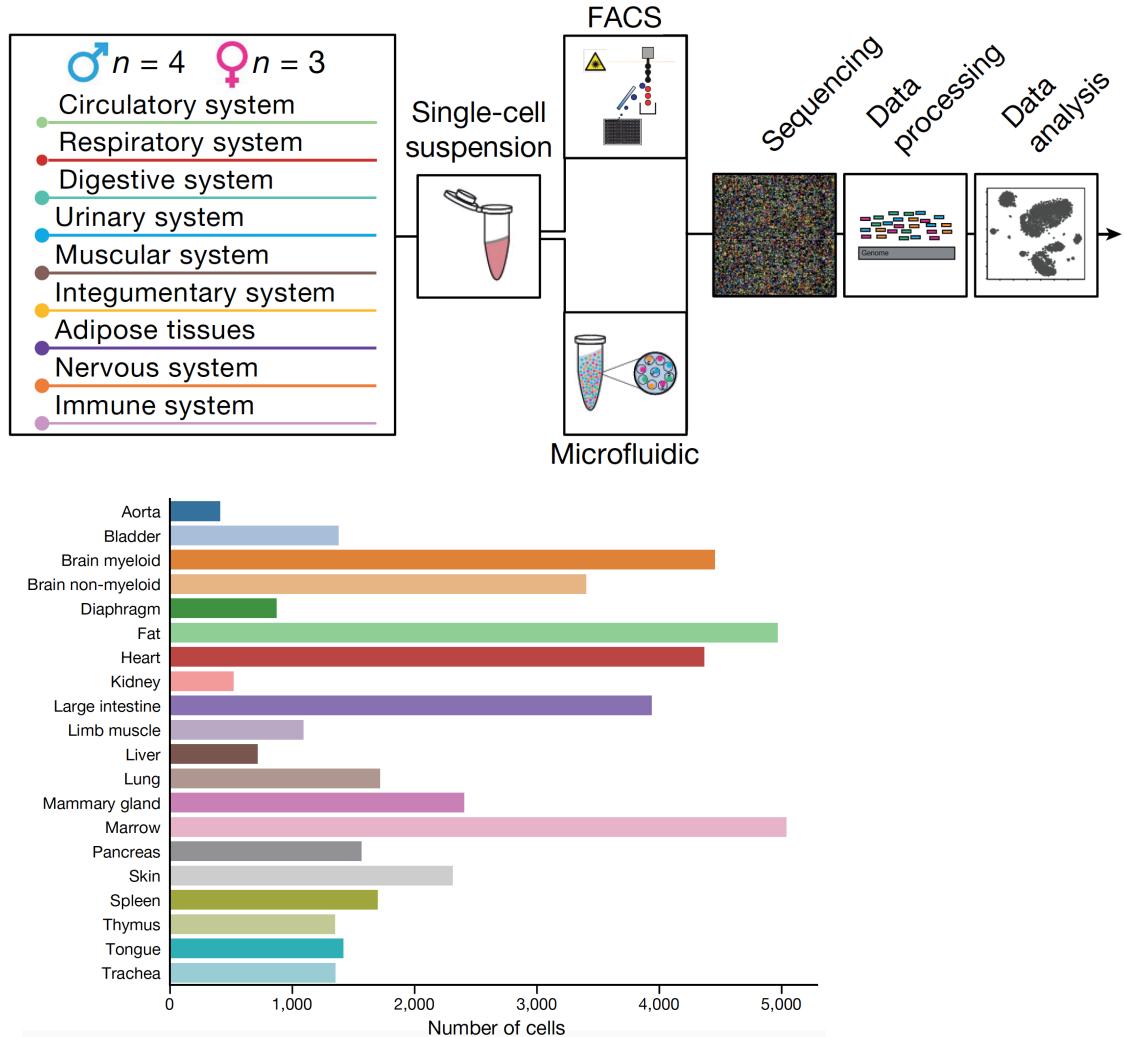


<https://doi.org/10.1038/s41586-018-0590-4>

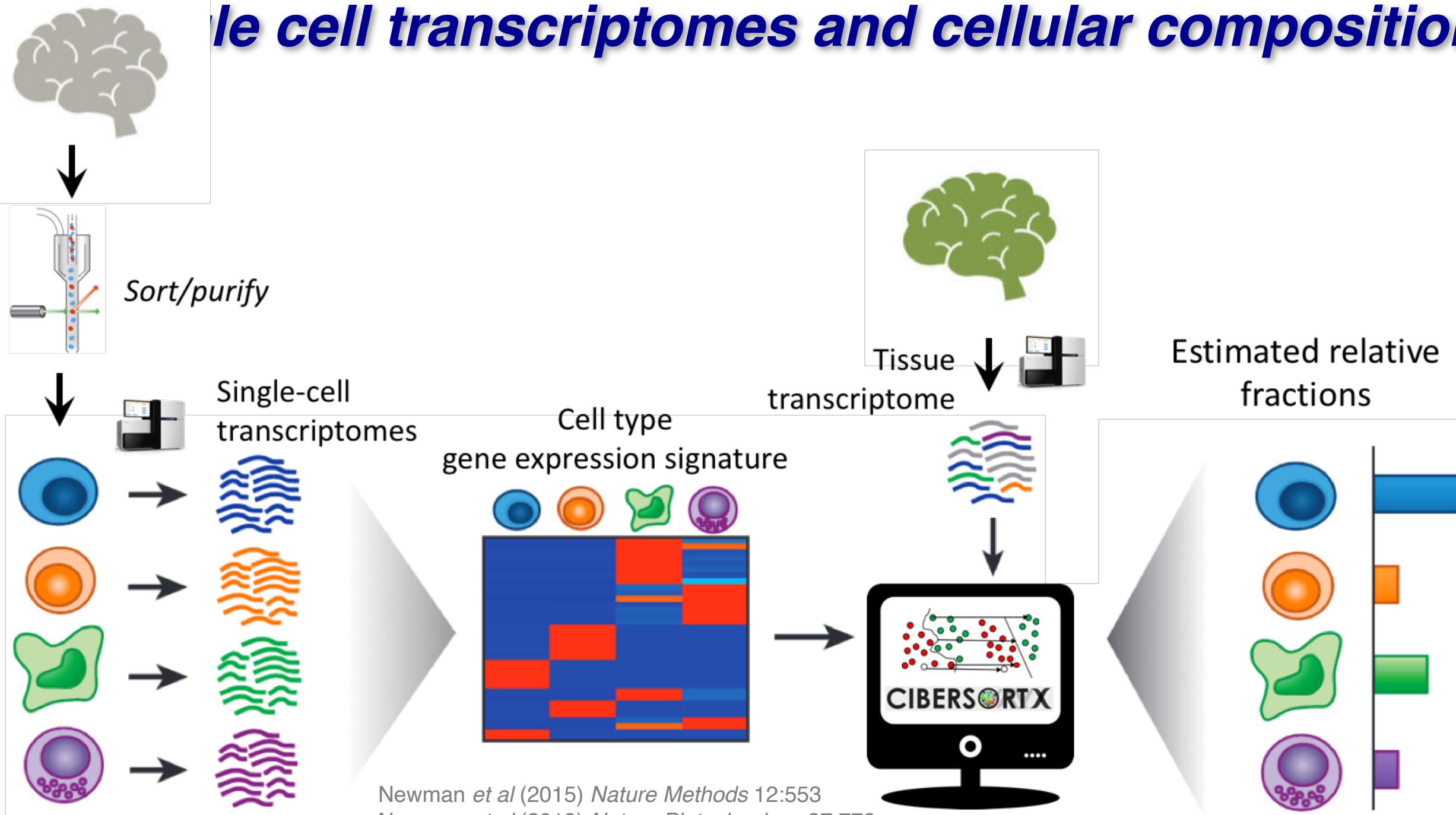
18 OCTOBER 2018 | VOL 562 | NATURE | 367

Single-cell transcriptomics of 20 mouse organs creates a *Tabula Muris*

The Tabula Muris Consortium*

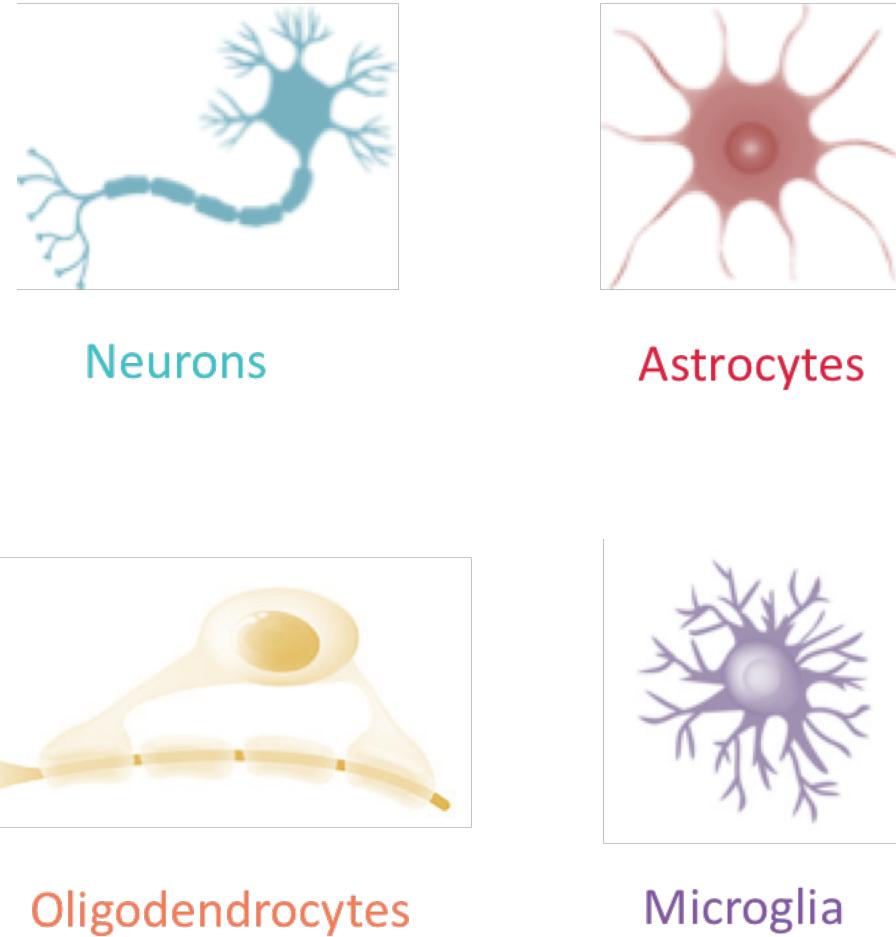
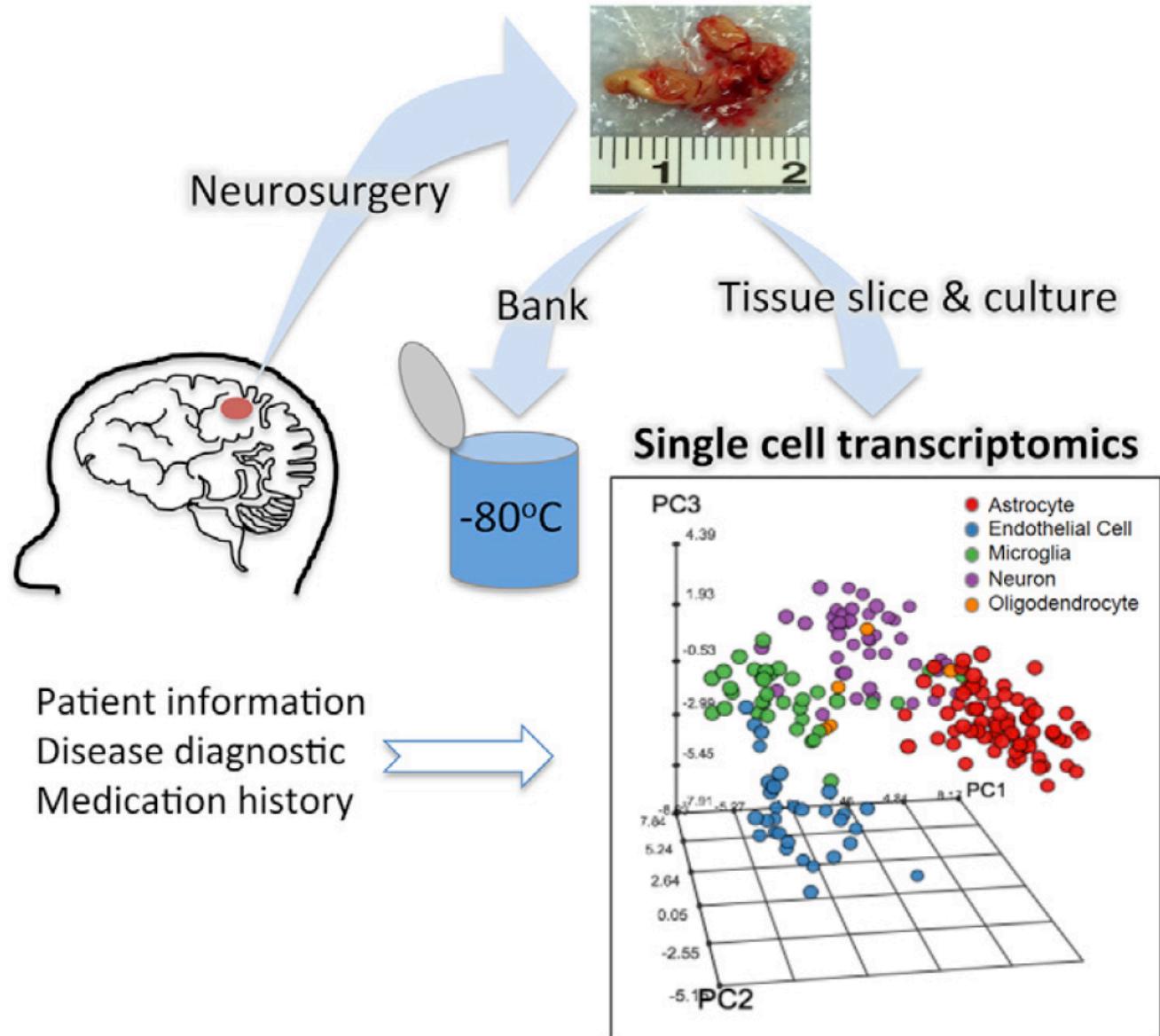


Single cell transcriptomes and cellular composition

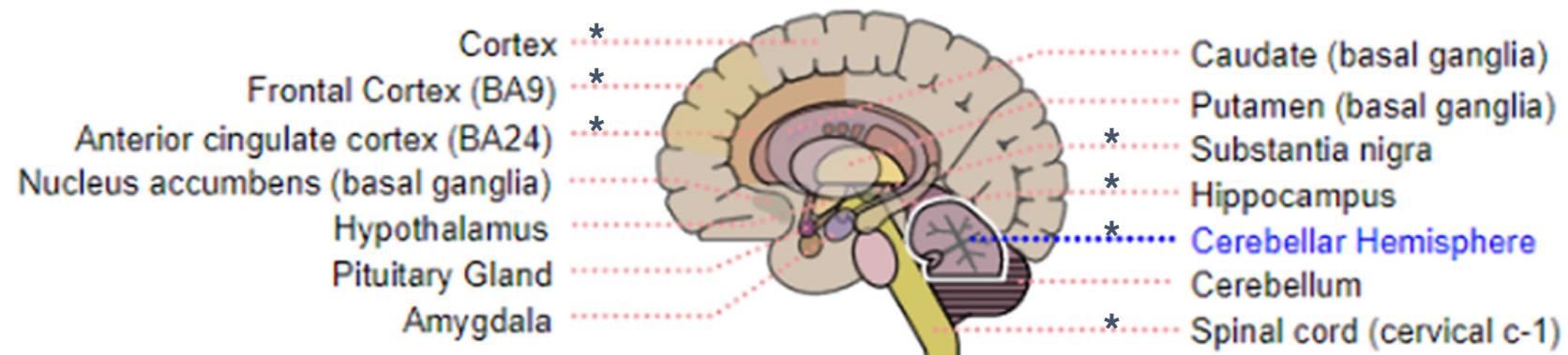


Human brain cell transcriptomes

Darmanis et al (2015) PNAS 112:7285
Spaethling et al (2017) Cell Reports 18:791



The ageing brain

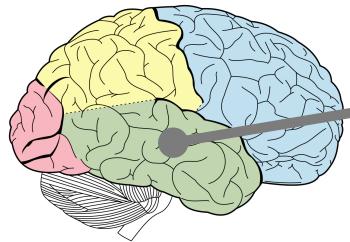


	Cortex		Frontal cortex		Substantia nigra		Anterior cingulate cortex		Hippocampus		Cerebellar hemisphere		Spinal cord	
	rho	p-value	rho	p-value	rho	p-value	rho	p-value	rho	p-value	rho	p-value	rho	p-value
	-0.26	0.00087	-0.19	0.031	-0.23	0.032	-0.22	0.016	-0.33	0.00016	0.045	0.6	-0.29	0.0053
	0.24	0.0023	0.2	0.024	0.17	0.1	0.2	0.031	0.28	0.0014	0.0013	0.99	0.16	0.12
	0.073	0.36	0.18	0.046	0.11	0.32	0.16	0.076	0.23	0.011	-0.31	0.00023	0.097	0.36
	0.11	0.18	0.082	0.36	0.091	0.4	0.25	0.0064	0.43	6.6E-07	0.03	0.73	0.3	0.004

positive correlation with age

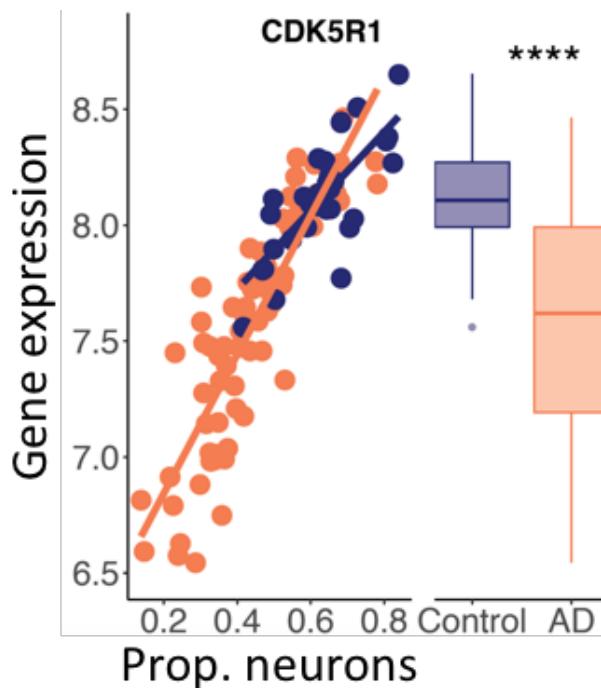
negative correlation with age

Loss of neurons in Alzheimer's disease (AD)



Temporal cortex RNA-seq
84 AD samples
54 Control samples

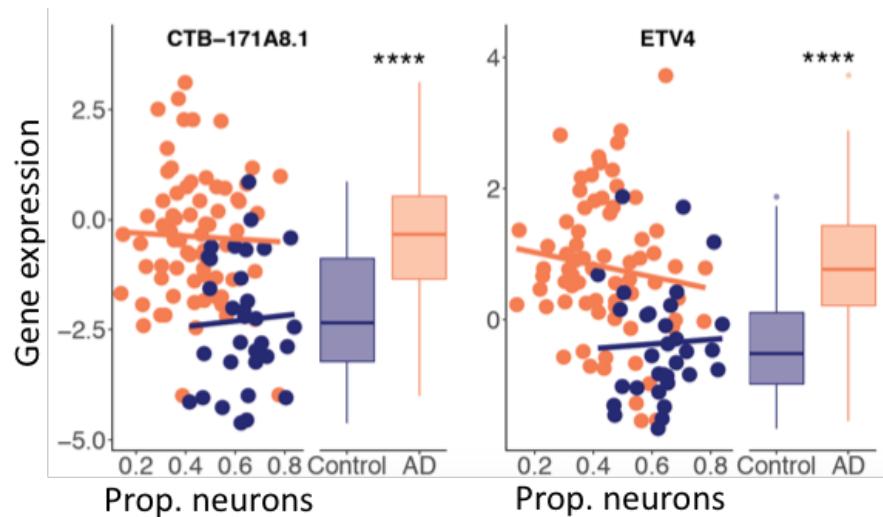
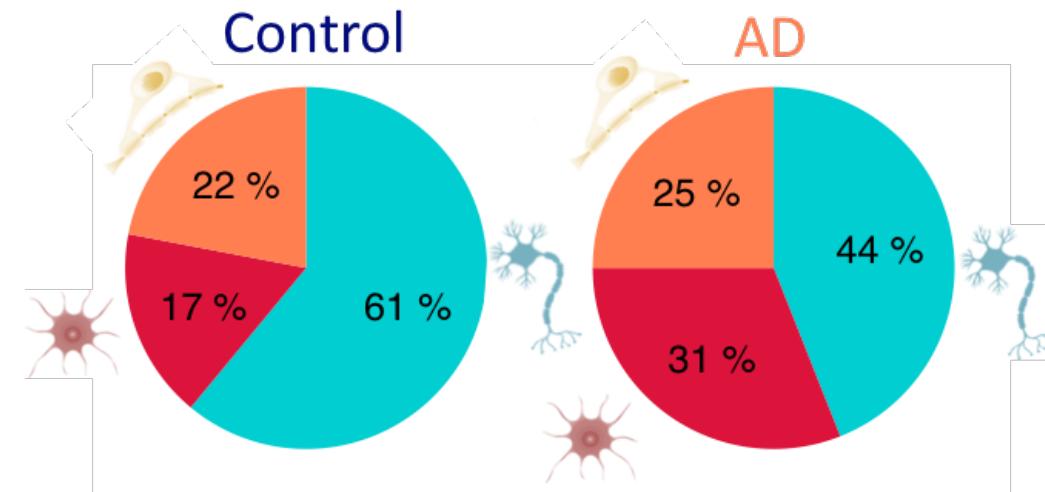
AMP-AD Knowledge
Portal
Pluta et al (2016) *Scientific Data* 3:160089



Mol Neurobiol
DOI 10.1007/s12035-016-0002-4

The miR-15/107 Family of microRNA Genes Regulates CDK5R1/p35 with Implications for Alzheimer's Disease Pathogenesis

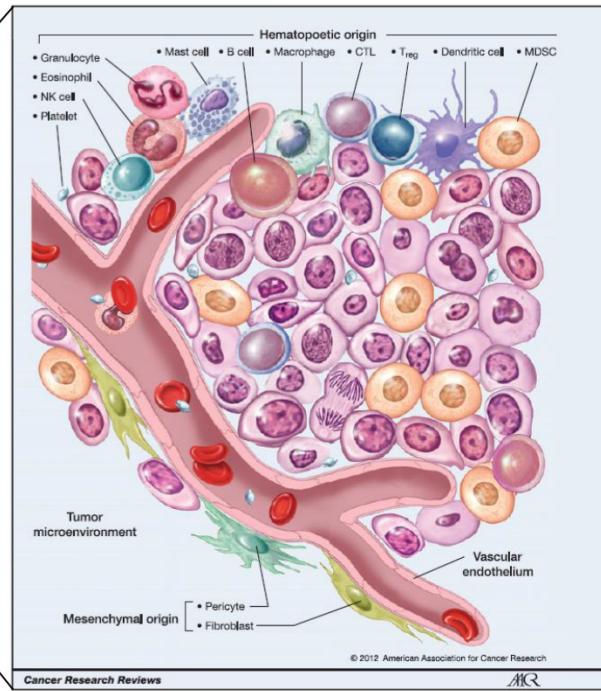
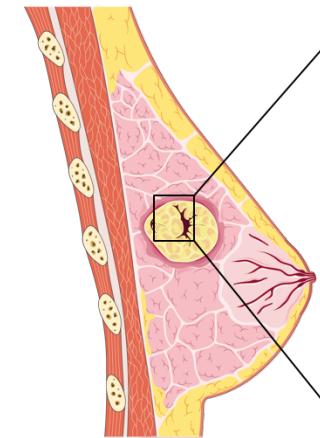
Silvia Moncini¹ · Marta Lunghi¹ · Alice Valmadre¹ · Margherita Grasso² ·
Valerio Del Vescovo² · Paola Riva¹ · Michela Alessandra Denti^{2,3} · Marco Venturin¹



Immune infiltration in breast cancer

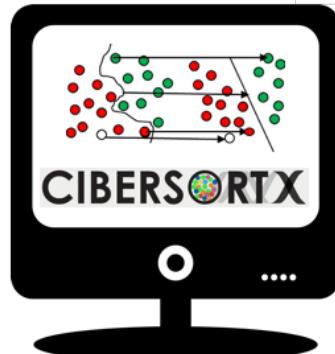
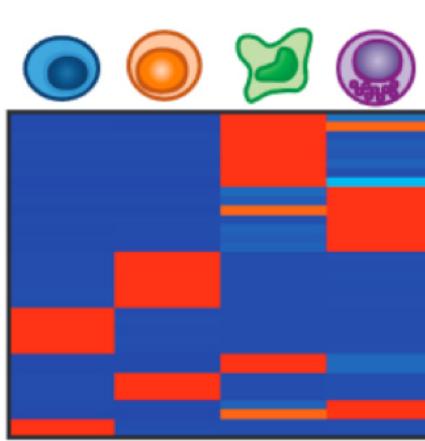


Marta



Kerkar and Restifo (2012) *Cancer Research*
72:3125

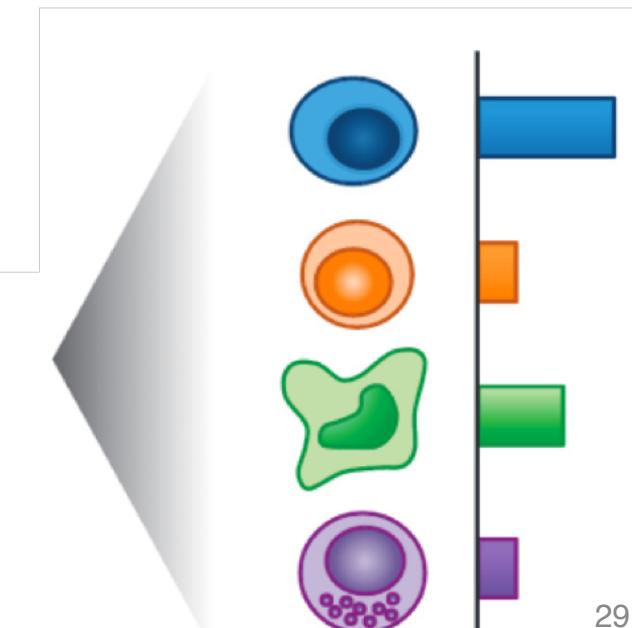
LM22:
leucocyte
(immune cell)
gene signature



Breast tumour
transcriptomes



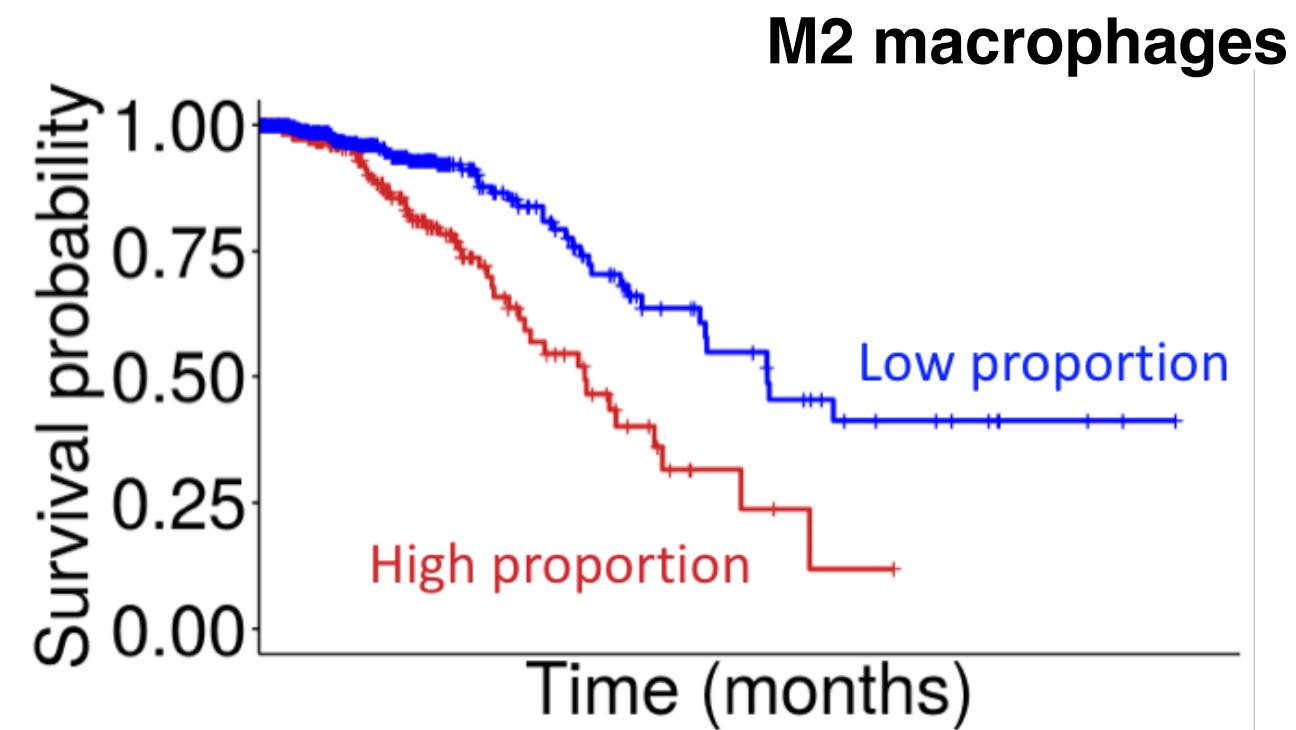
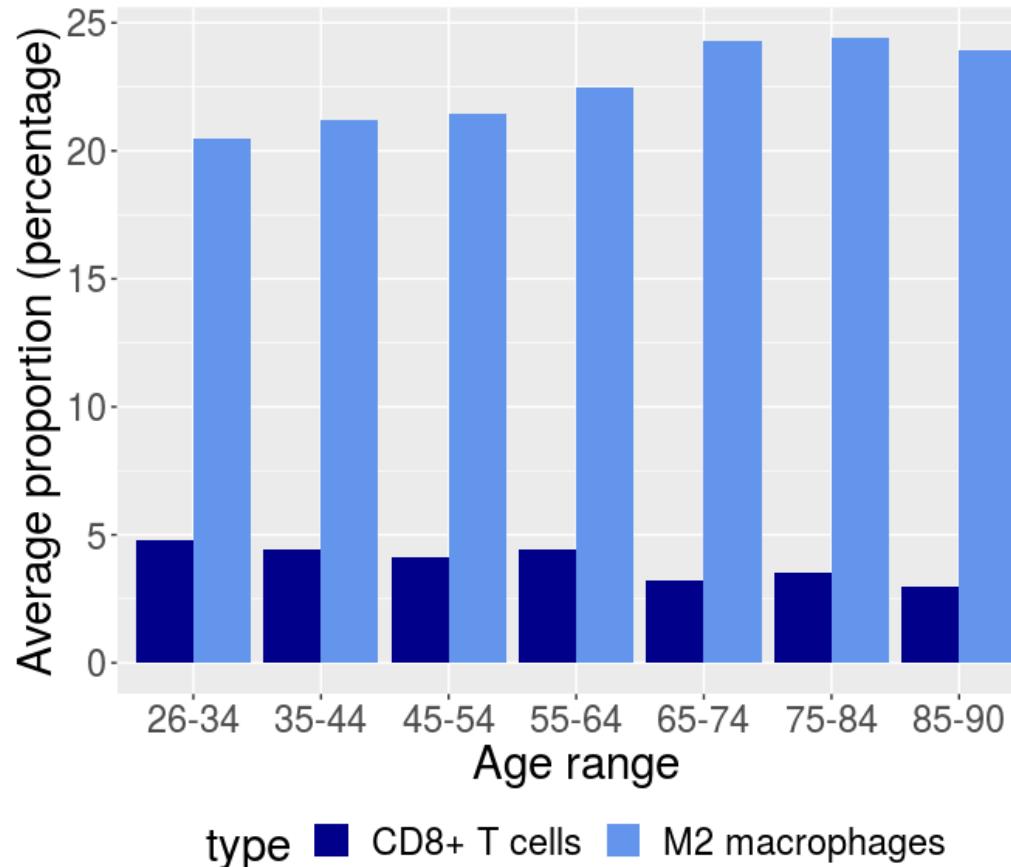
Estimated relative
fractions



Immune infiltration and age in breast cancer



- CD8+ T cells: toxic to other cells, anti-tumoural
- M2 macrophages: tissue repair, pro-tumourigenic

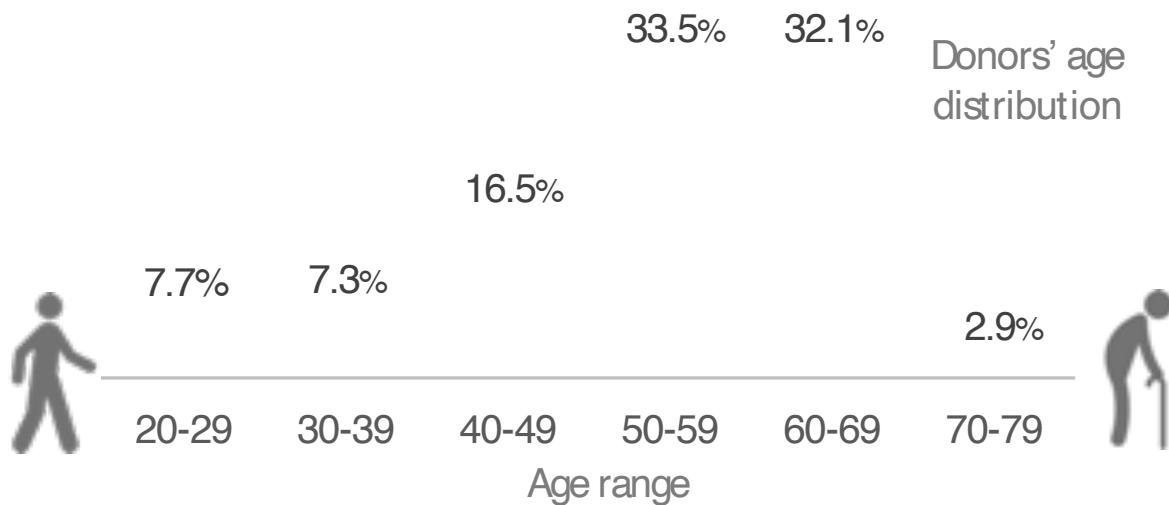


type ■ CD8+ T cells ■ M2 macrophages

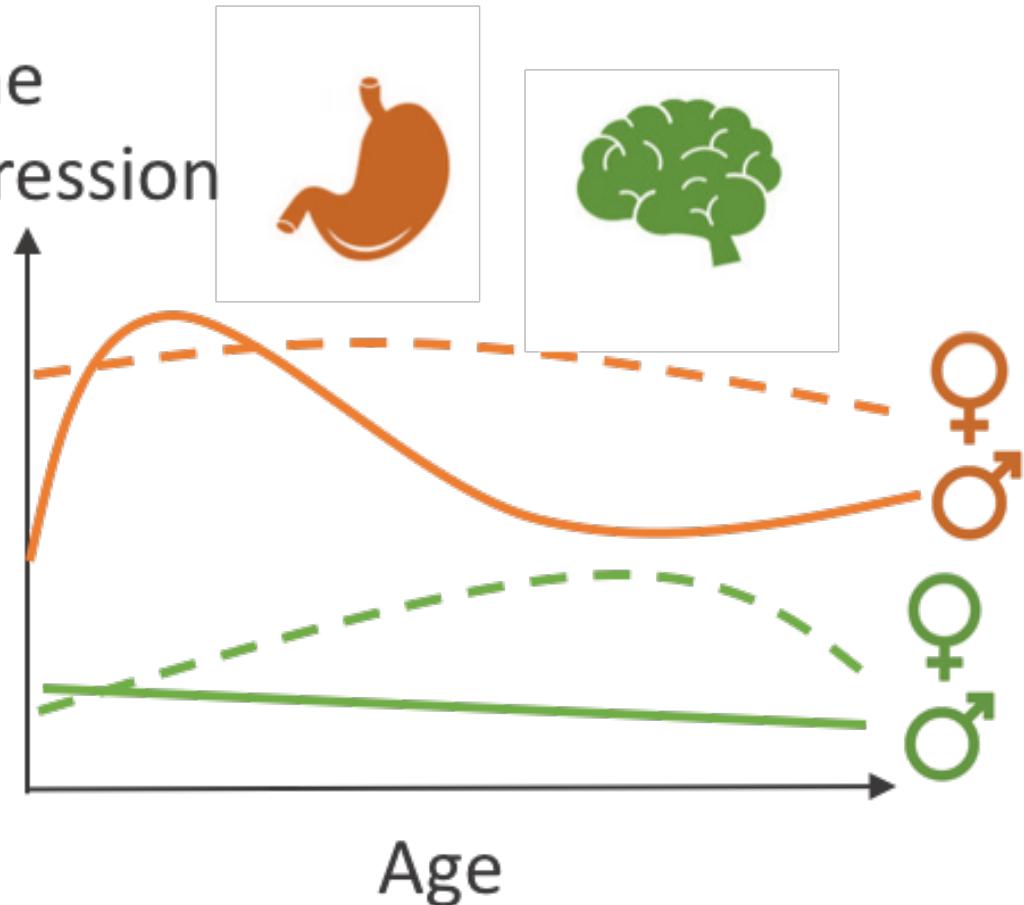
How gene expression varies with age in human tissues



Arthur

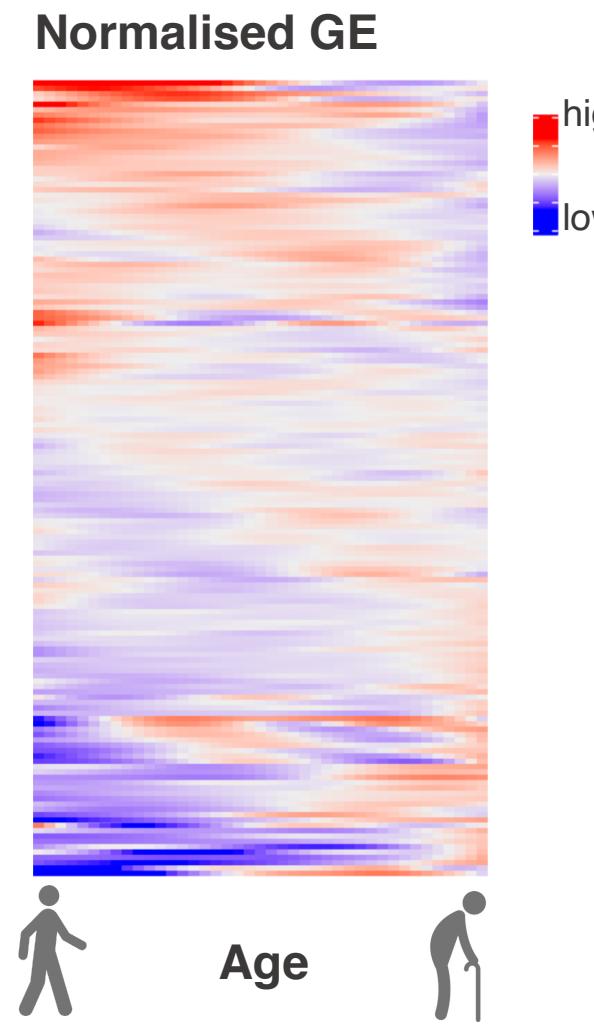
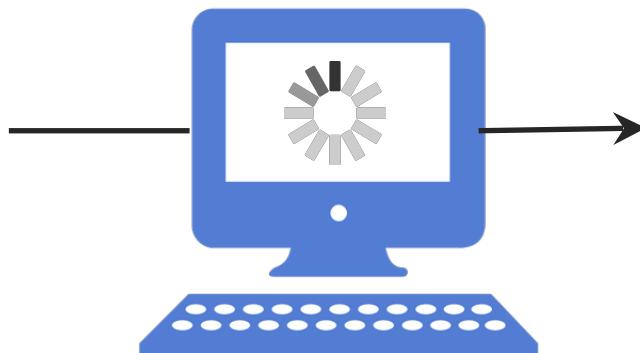
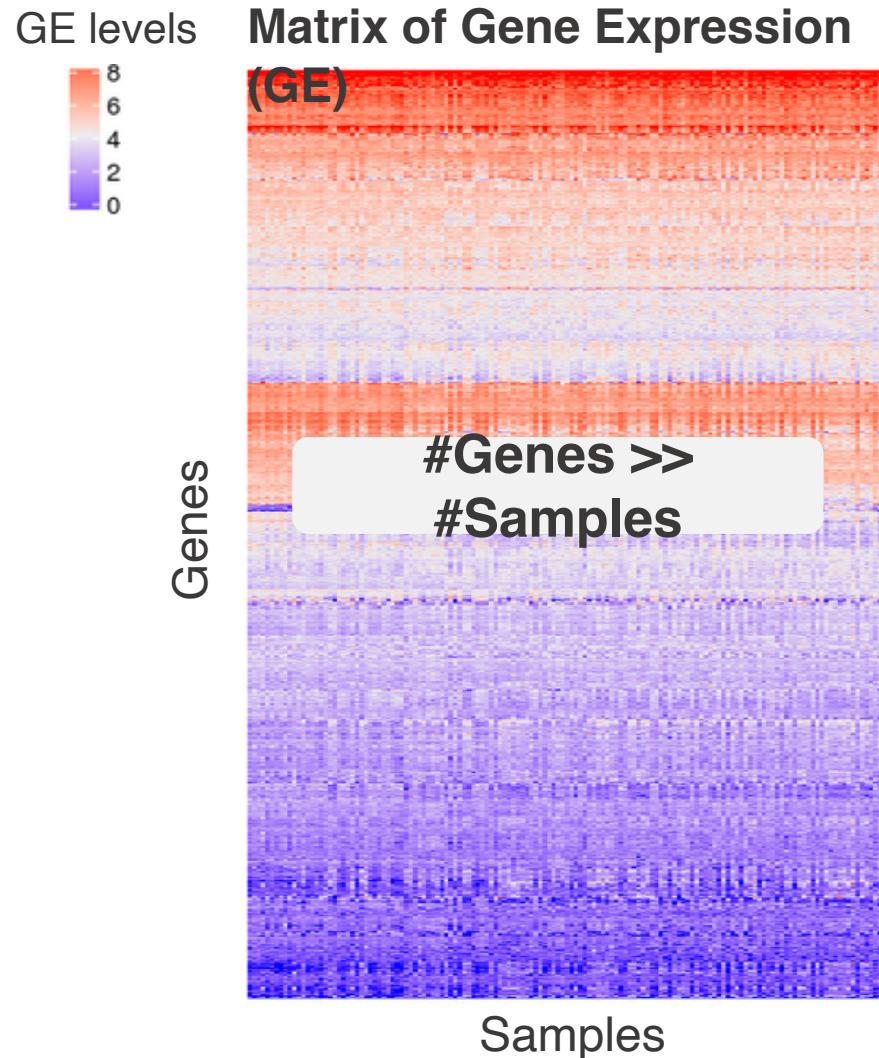


Gene Expression

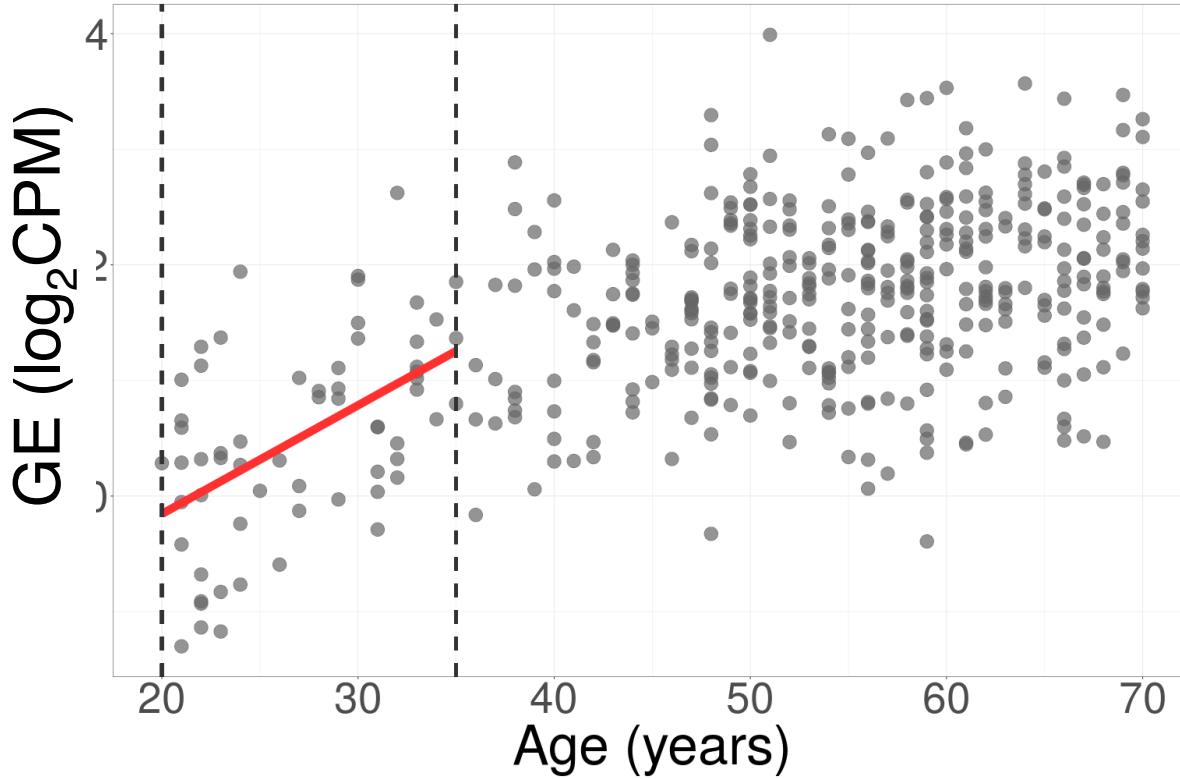




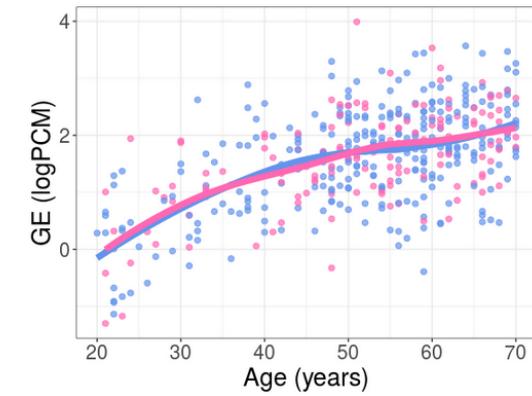
For each tissue/gender combination:



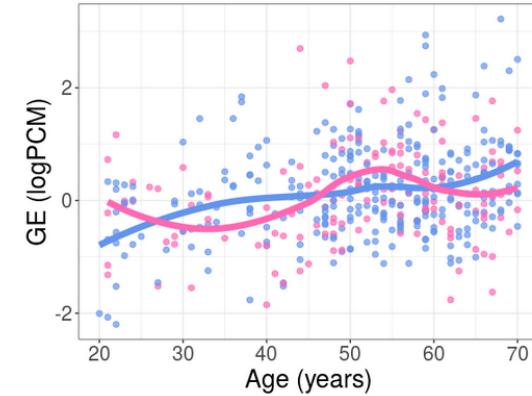
For each tissue/gender combination:



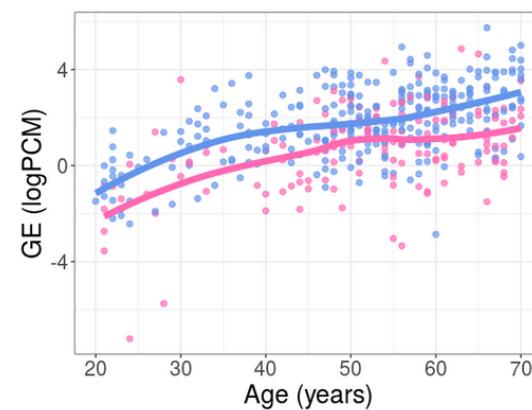
Gender differences?



None

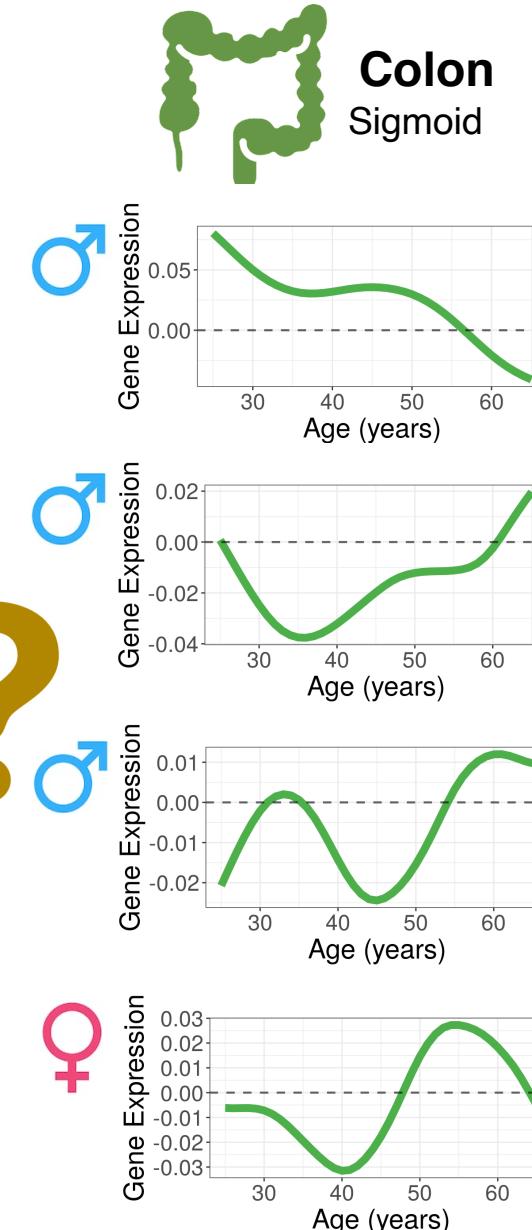
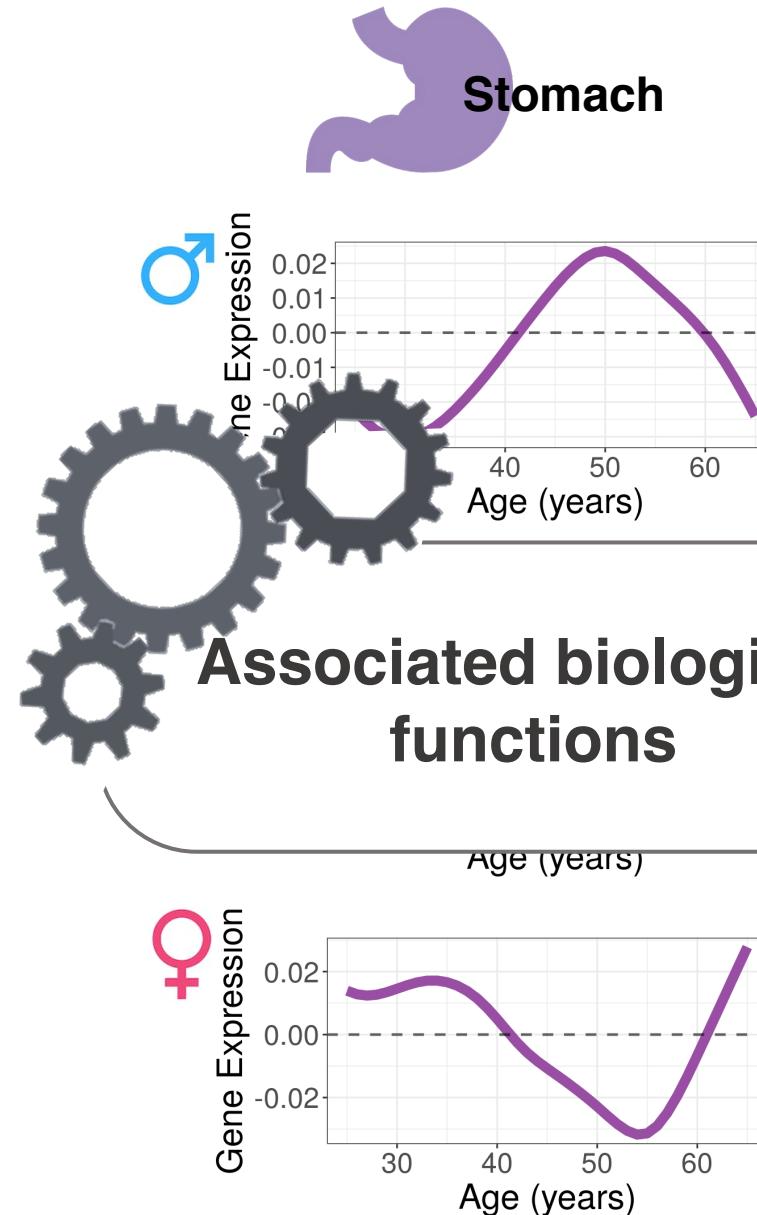
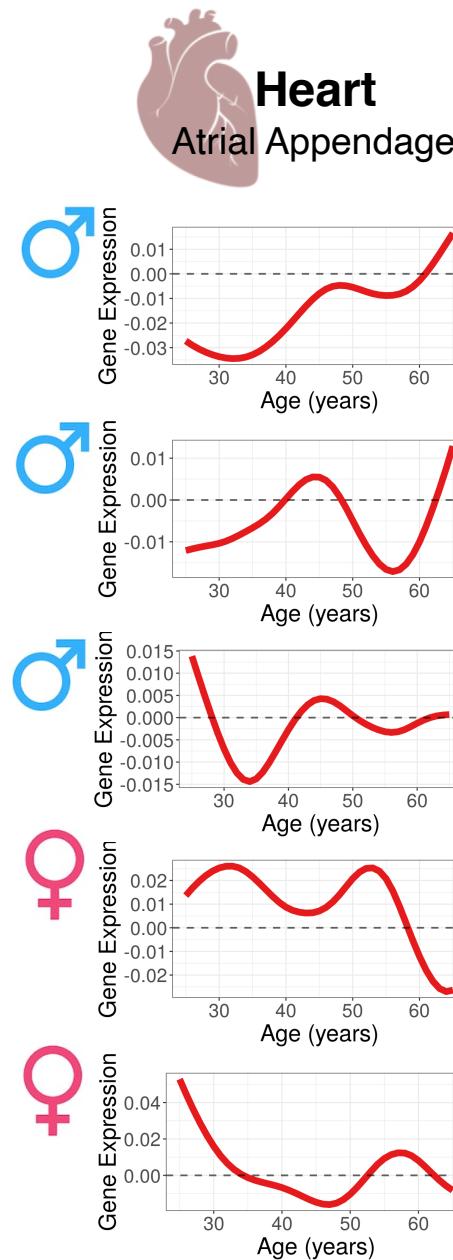


Trend



Magnitude

Tissue/gender-specific ageing



Associated biological functions

Tissue-specific ageing



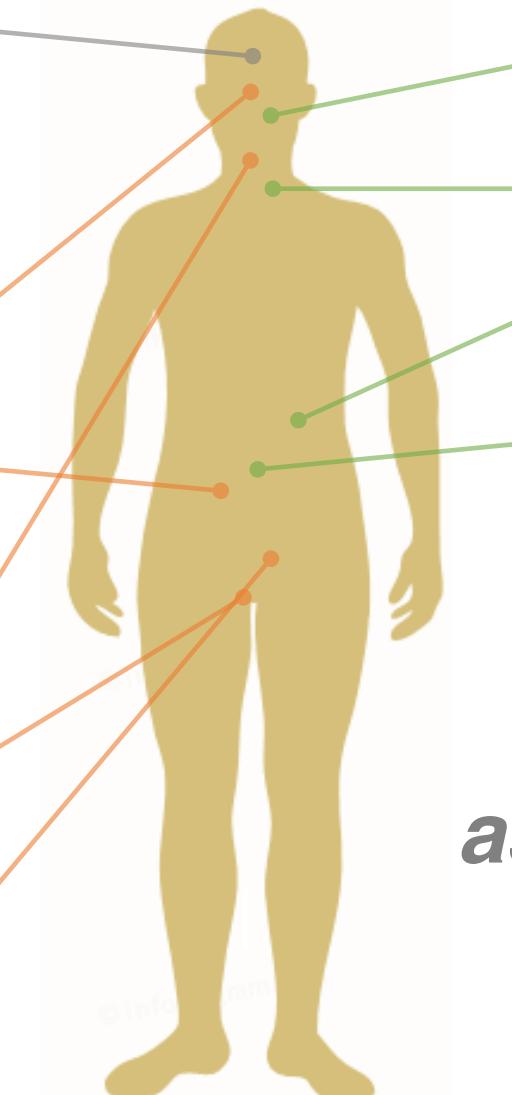
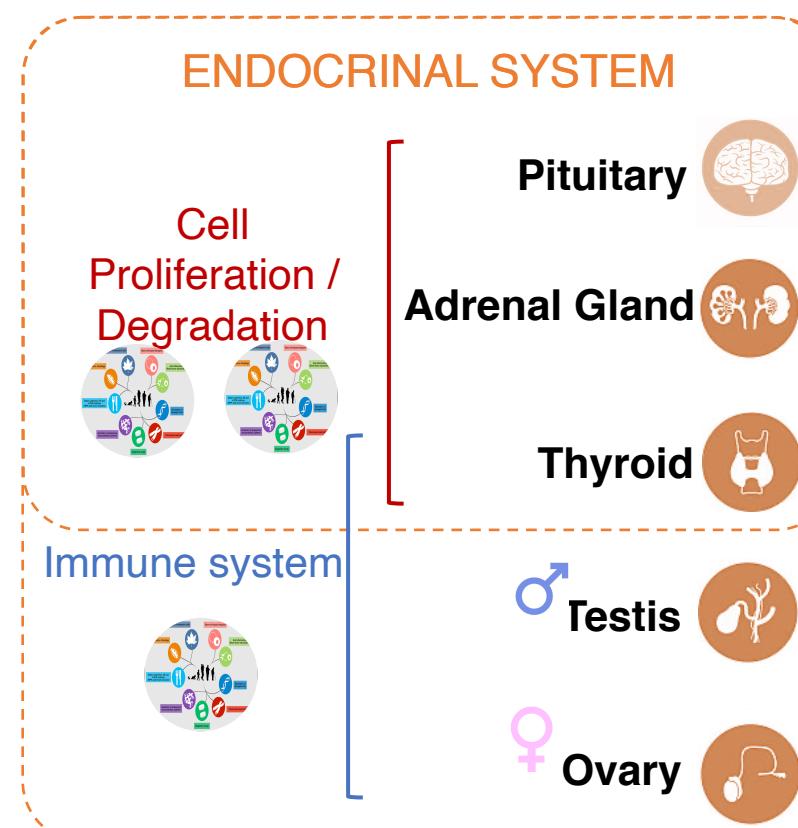
Immune Response
Cell Proliferation
and Maintenance

Brain

Depleted in driver genes
Enriched in driver genes



Arthur



METABOLIC FUNCTIONS

- Minor Salivary Gland
- Oesophagus
- Colon
- Stomach
- Adipose

Cell Proliferation
and Maintenance



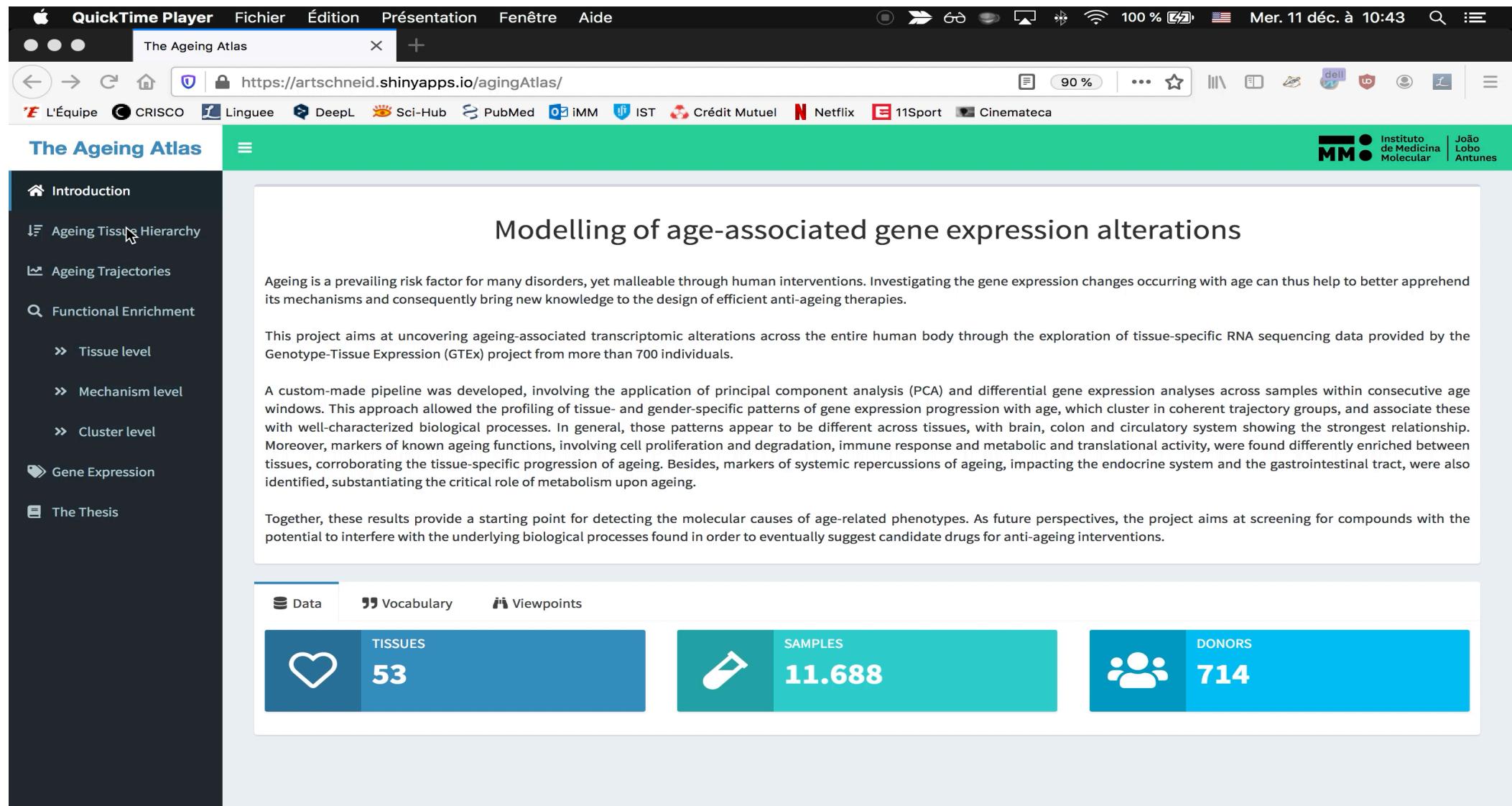
*asynchronous inter-organ
progression of human
ageing*

Future: the Ageing Atlas

<https://artschneid.shinyapps.io/agingAtlas/>



Arthur



The screenshot shows a web browser window titled "The Ageing Atlas". The URL in the address bar is <https://artschneid.shinyapps.io/agingAtlas/>. The page content is as follows:

Modelling of age-associated gene expression alterations

Ageing is a prevailing risk factor for many disorders, yet malleable through human interventions. Investigating the gene expression changes occurring with age can thus help to better apprehend its mechanisms and consequently bring new knowledge to the design of efficient anti-ageing therapies.

This project aims at uncovering ageing-associated transcriptomic alterations across the entire human body through the exploration of tissue-specific RNA sequencing data provided by the Genotype-Tissue Expression (GTEx) project from more than 700 individuals.

A custom-made pipeline was developed, involving the application of principal component analysis (PCA) and differential gene expression analyses across samples within consecutive age windows. This approach allowed the profiling of tissue- and gender-specific patterns of gene expression progression with age, which cluster in coherent trajectory groups, and associate these with well-characterized biological processes. In general, those patterns appear to be different across tissues, with brain, colon and circulatory system showing the strongest relationship. Moreover, markers of known ageing functions, involving cell proliferation and degradation, immune response and metabolic and translational activity, were found differently enriched between tissues, corroborating the tissue-specific progression of ageing. Besides, markers of systemic repercussions of ageing, impacting the endocrine system and the gastrointestinal tract, were also identified, substantiating the critical role of metabolism upon ageing.

Together, these results provide a starting point for detecting the molecular causes of age-related phenotypes. As future perspectives, the project aims at screening for compounds with the potential to interfere with the underlying biological processes found in order to eventually suggest candidate drugs for anti-ageing interventions.

Below the text, there are three cards:

- Data**: TISSUES 53
- Vocabulary**: SAMPLES 11.688
- Viewpoints**: DONORS 714

Future: translation

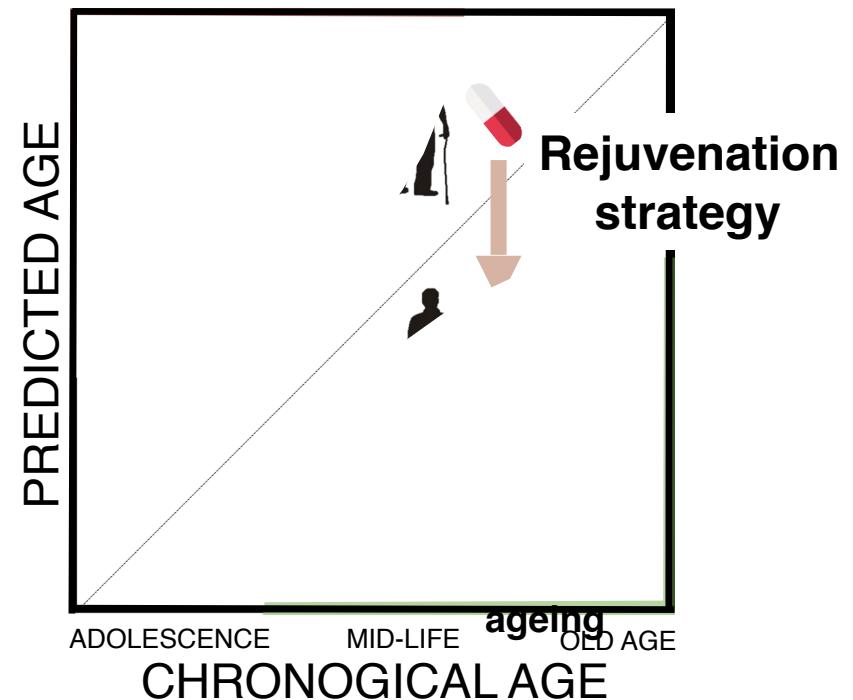


Arthur

1. Characterise causal **genetic** or **cellular composition** alterations responsible for the profiled ageing patterns

Pinpoint **pharmacological** interventions with potential to revert them

2. Develop machine learning model to **predict** the donor's age from transcriptomic data



3. Why biology needs more data scientists

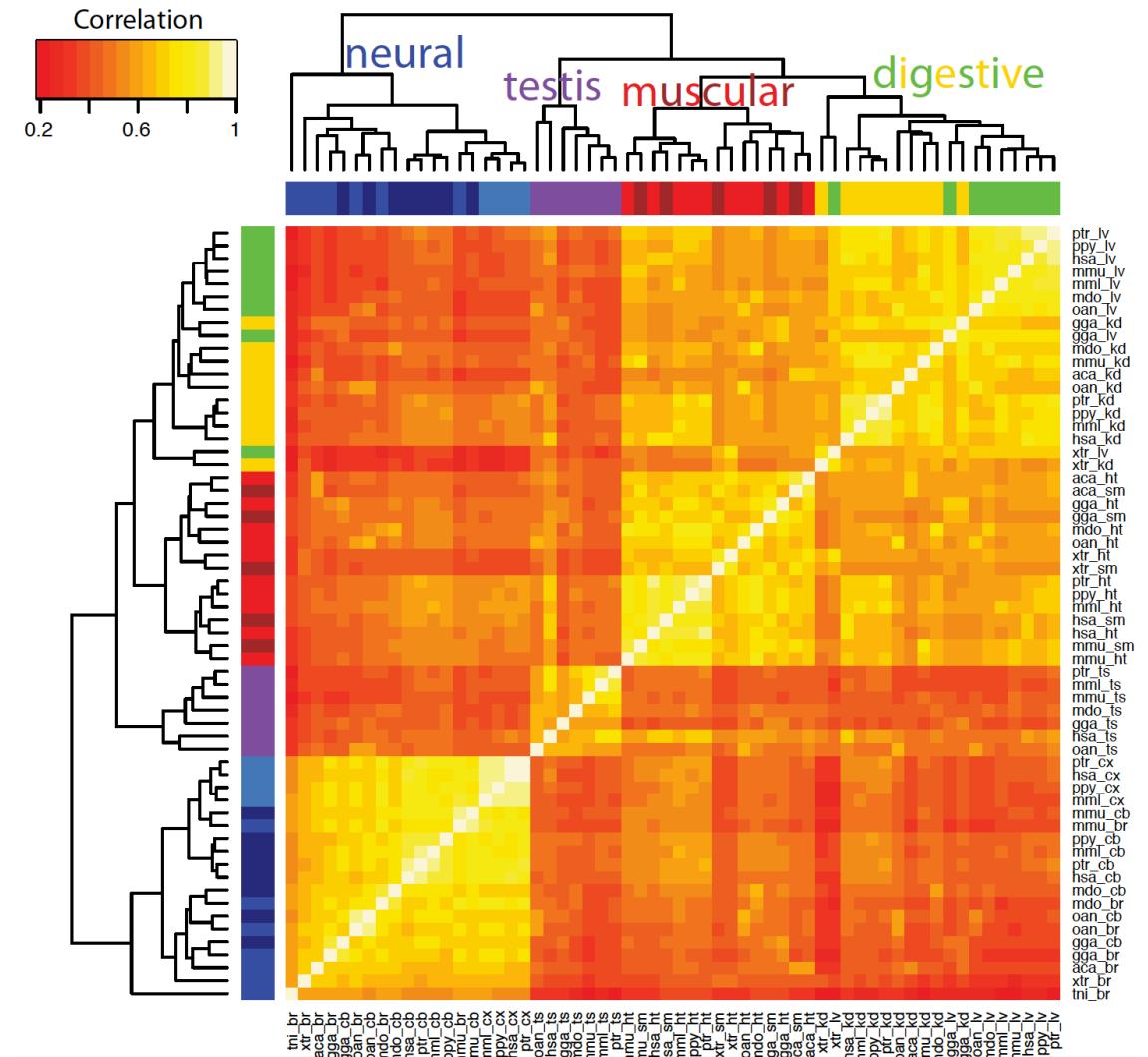
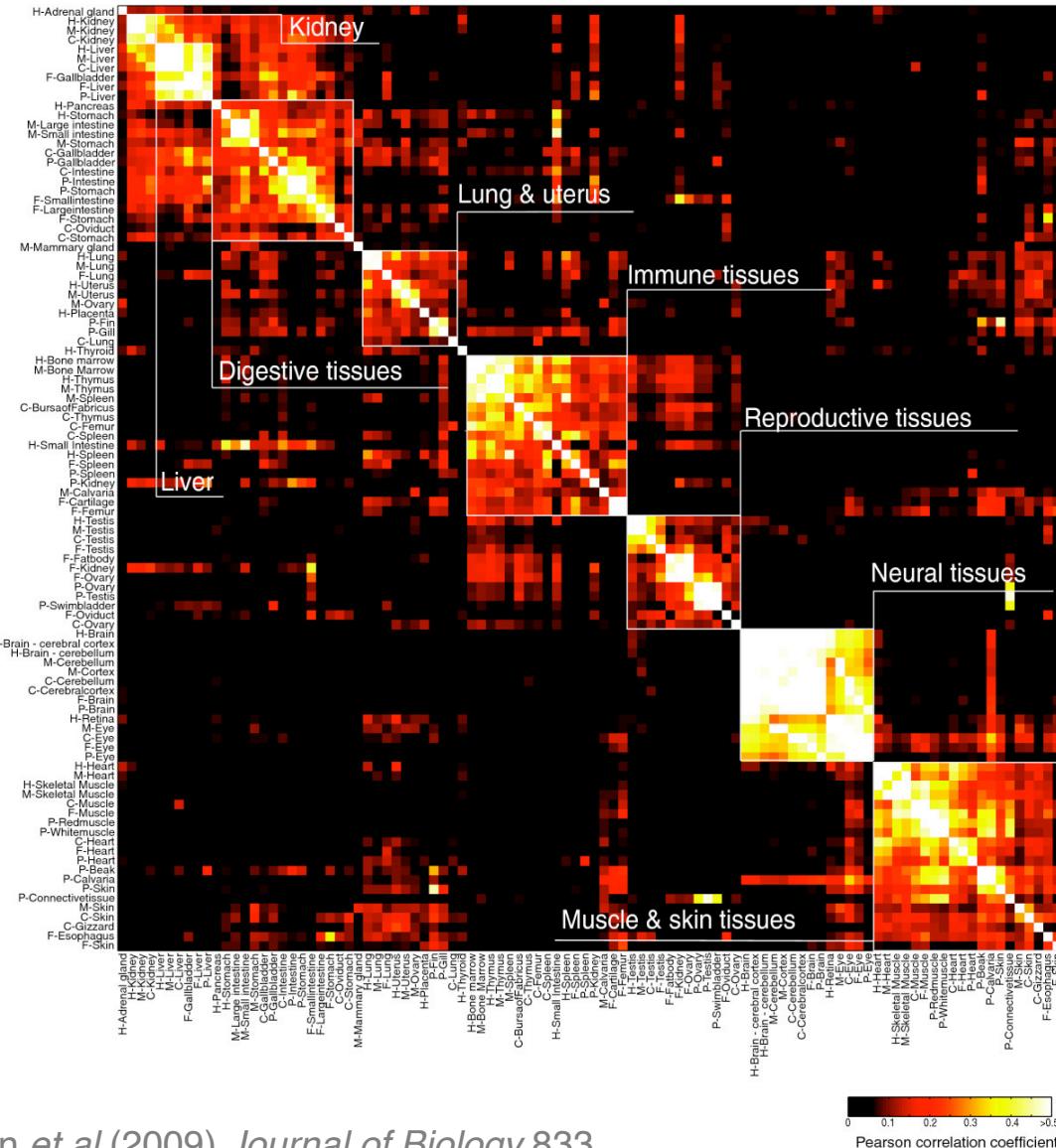


Tissue-specific gene expression is conserved in vertebrates

20 tissues/organs, 5 species, microarrays

vertebrates

7 tissues/organs, 11 species, RNA-seq

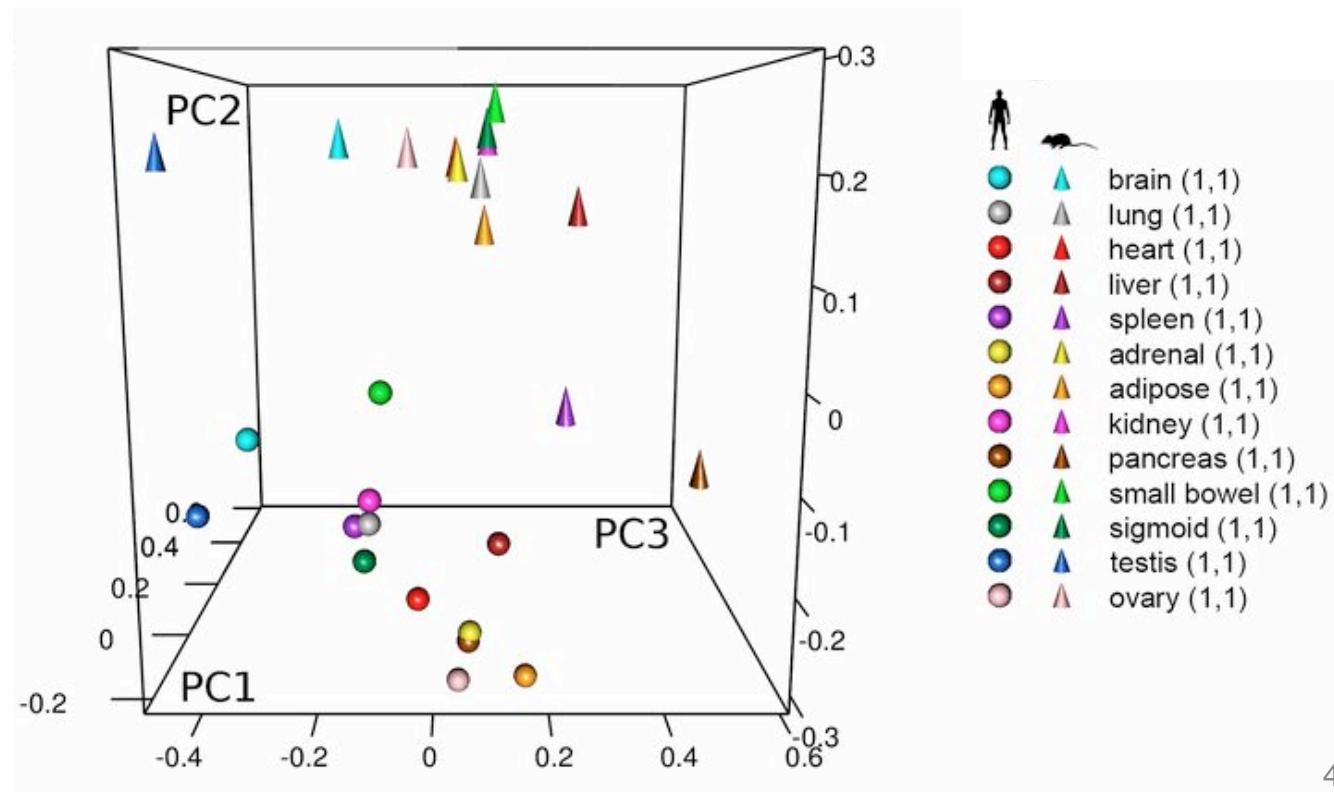
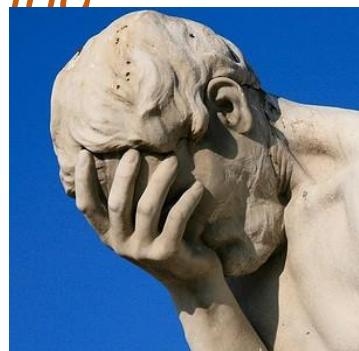


Gene Expression Is More Similar Among Tissues Within a Species Than Between Corresponding Tissues of the Two Species

Comparison of the transcriptional landscapes between human and mouse tissues

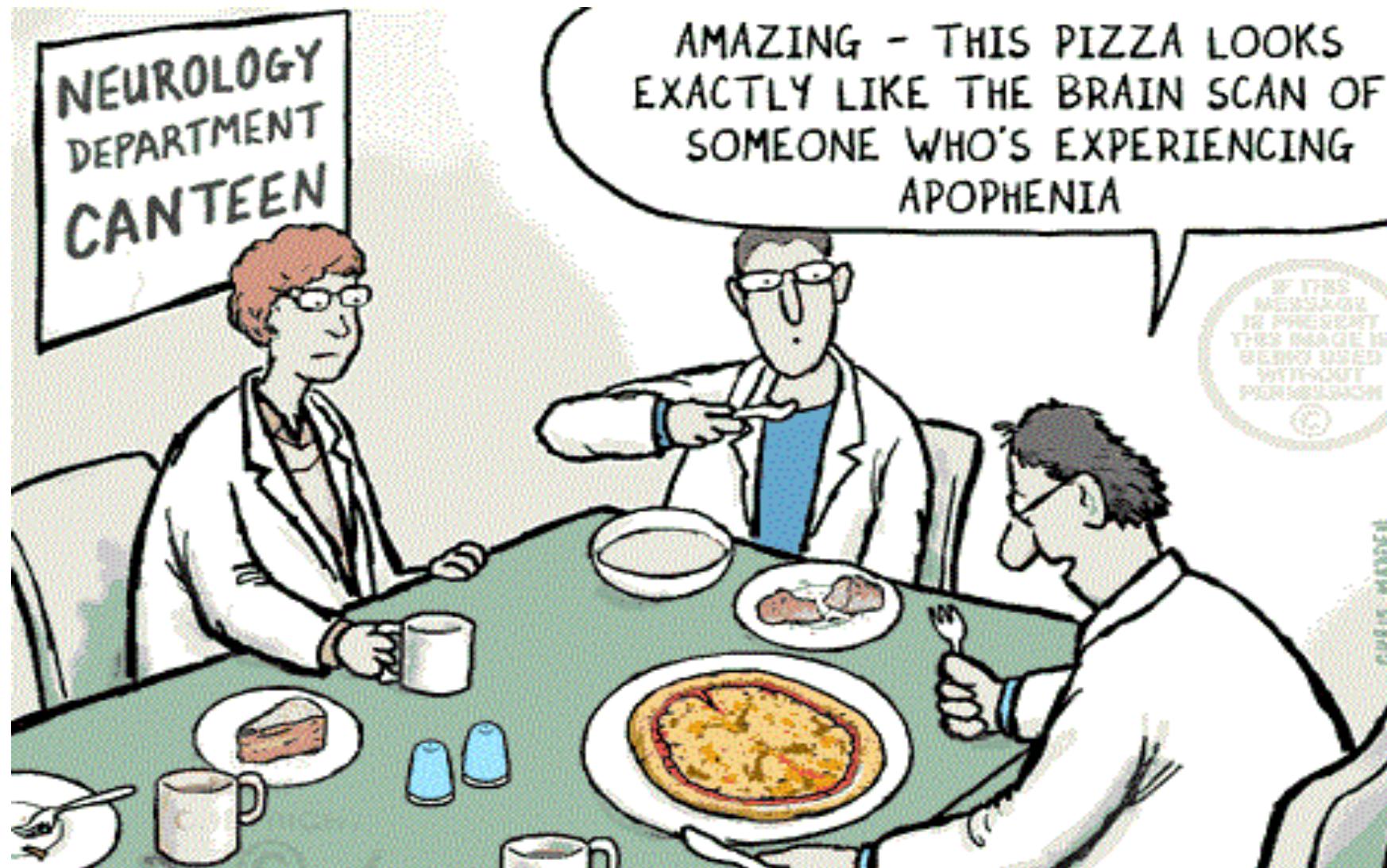
17224–17229 | PNAS | December 2, 2014 | vol. 111 | no. 48

analysis of only the 13 paired samples processed under one experimental protocol yielded the same species-specific clustering



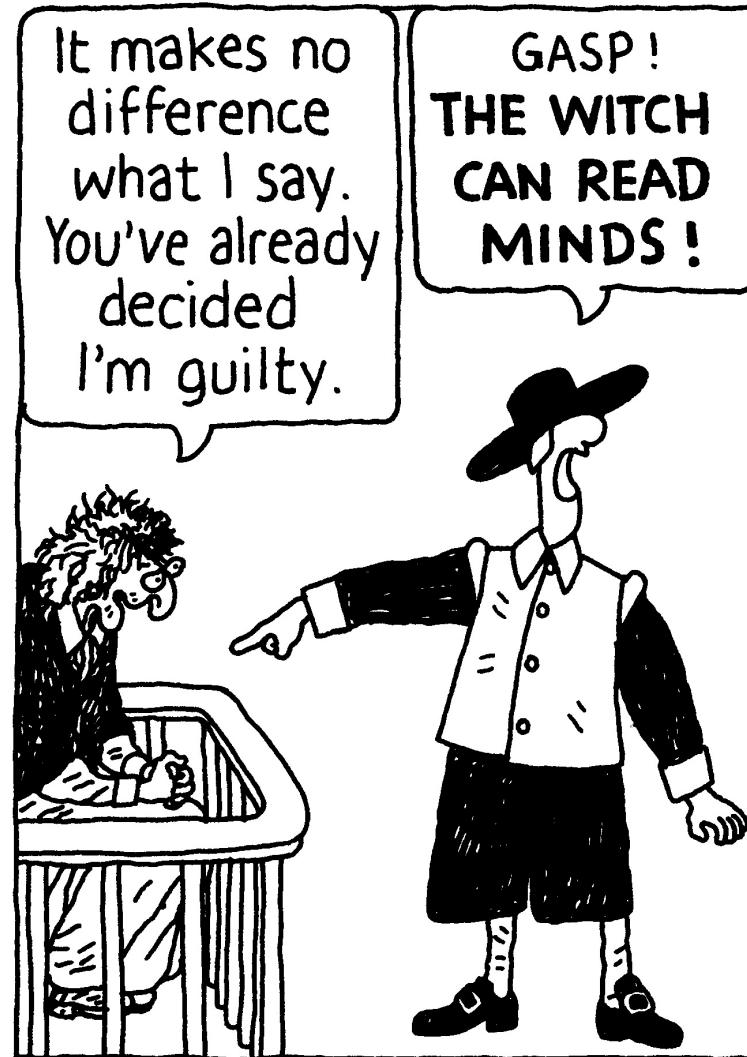
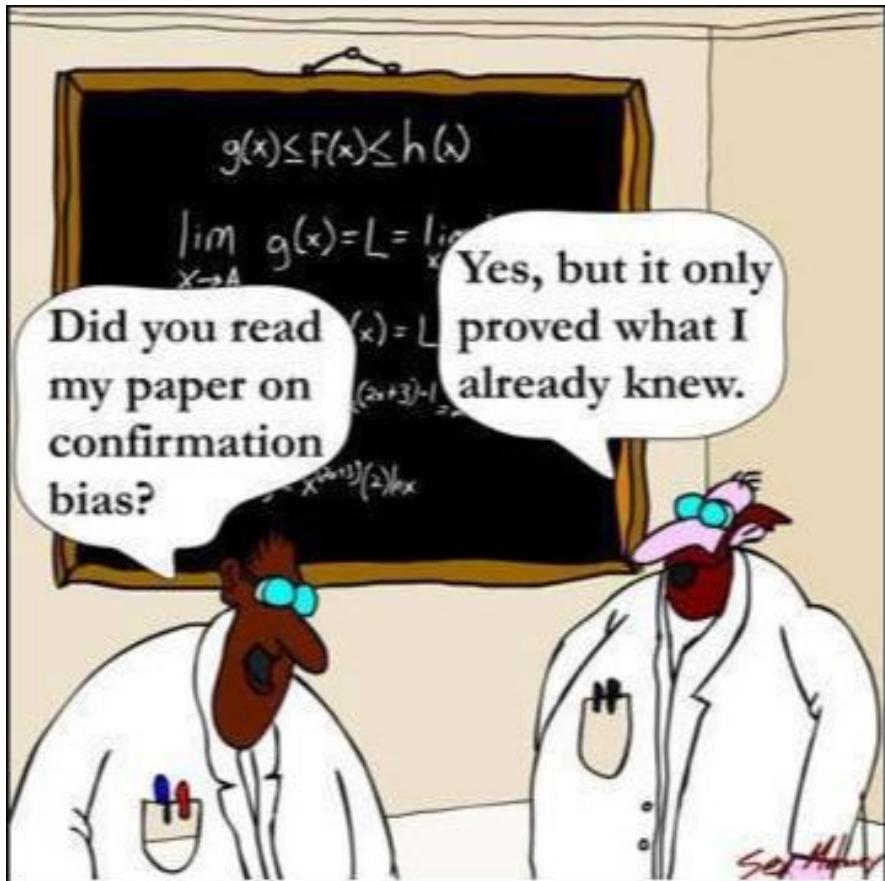
Apophenia

human tendency to perceive meaningful patterns within random data



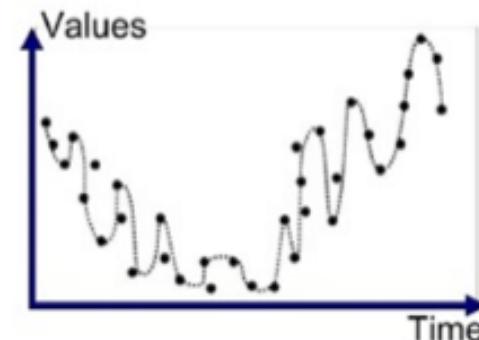
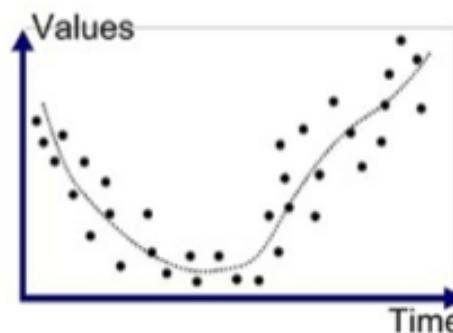
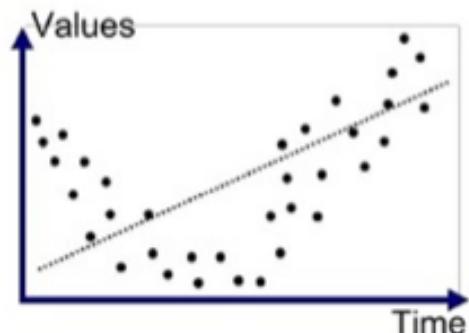
Confirmation bias

tendency to interpret new evidence as confirmation of one's existing beliefs or theories

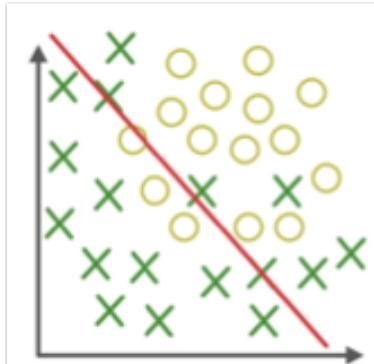


Overfitting

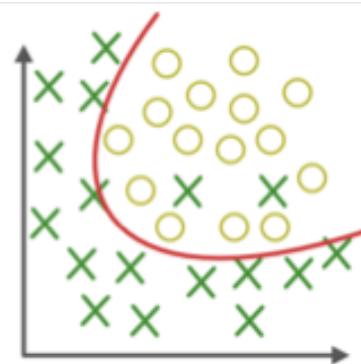
the statistical model describes random error or noise instead of the underlying relationship



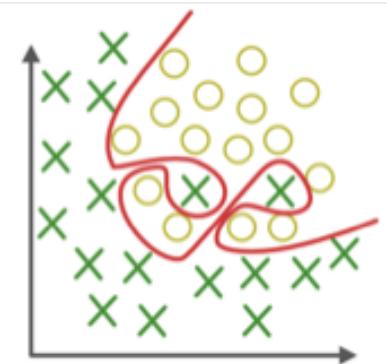
Underfitting



Properly fitting



Overfitting



An approximate answer to the right problem is worth a good deal more than an exact answer to an approximate problem.
John Tukey (1915-2000) - Statistician



Distraction

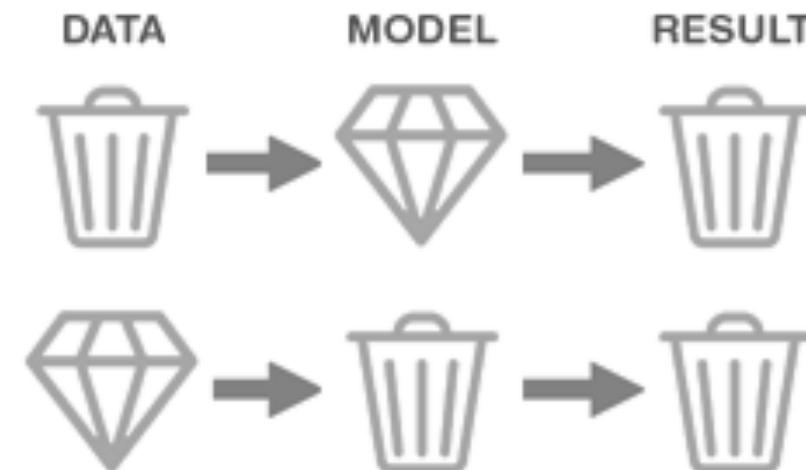
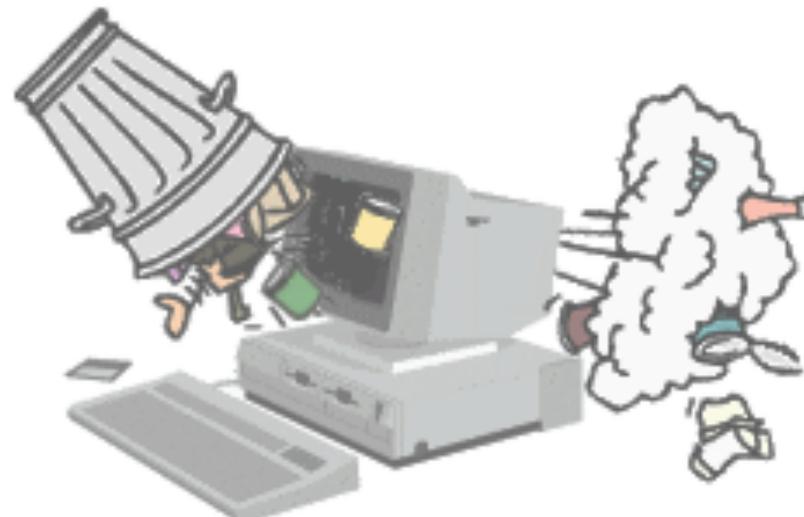
too much interesting information

(one always finds exciting biology, not necessarily related with the original question)

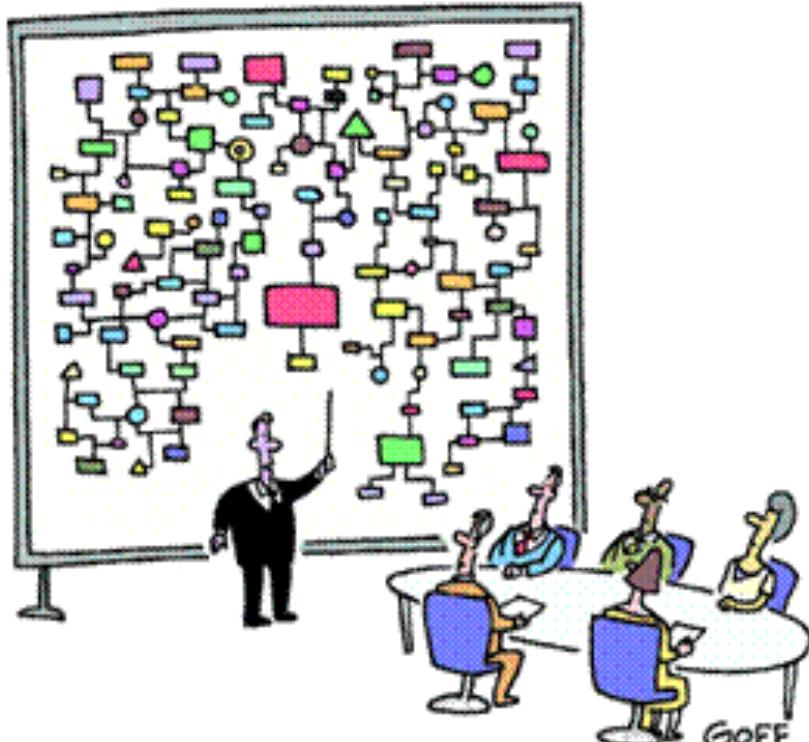


Complacency & illiteracy – the black box

- Not bothering with statistical concepts and assumptions
- Looking for easy **recipes**, blindly judging numbers (p-values) without “getting intimate” with data (e.g. plotting)
- Evaluation of fairness of method based on outcome



Misestimate of biological systems' complexity



"And that's why we need a computer."



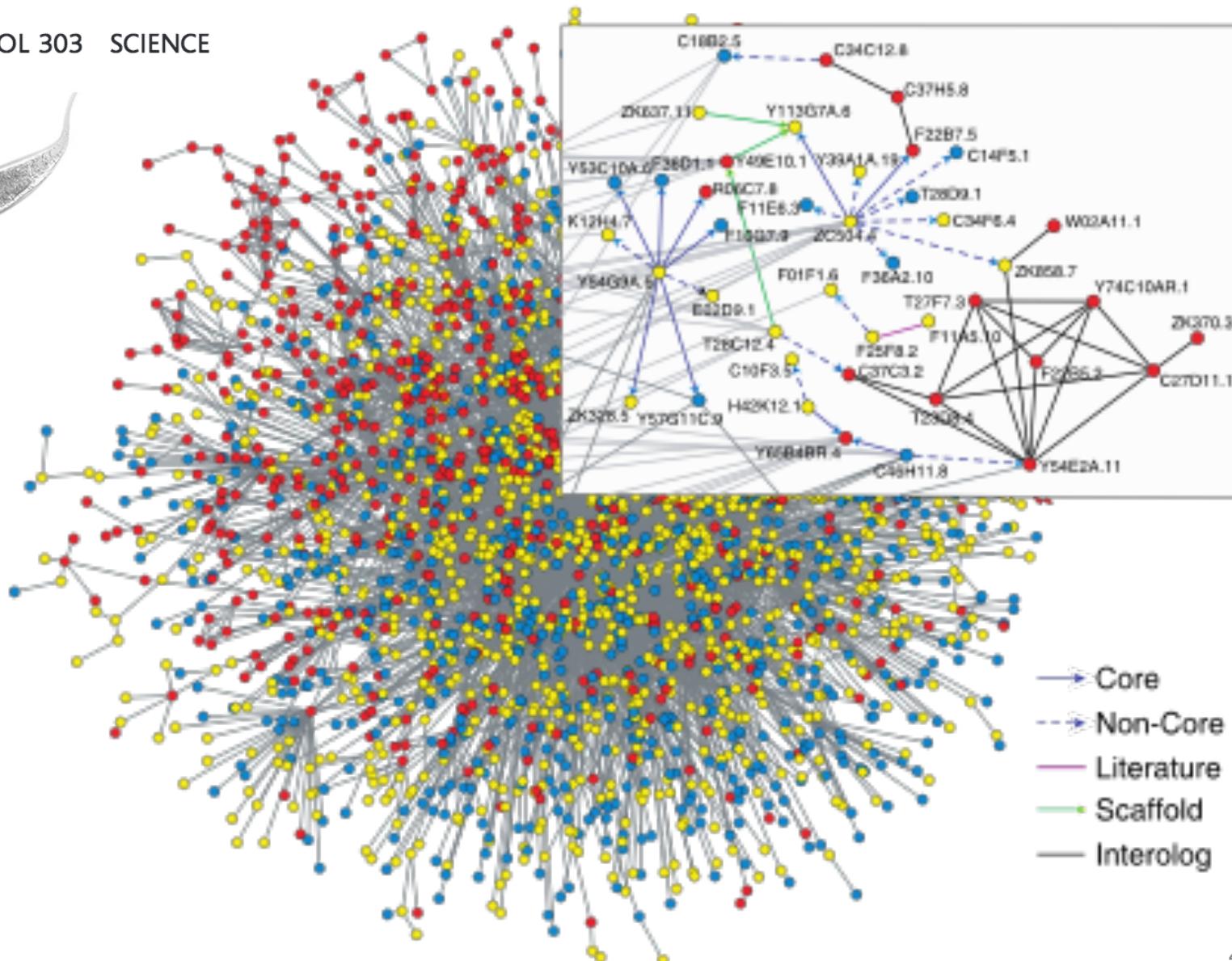
A Map of the Interactome Network of the Metazoan *C. elegans*

540

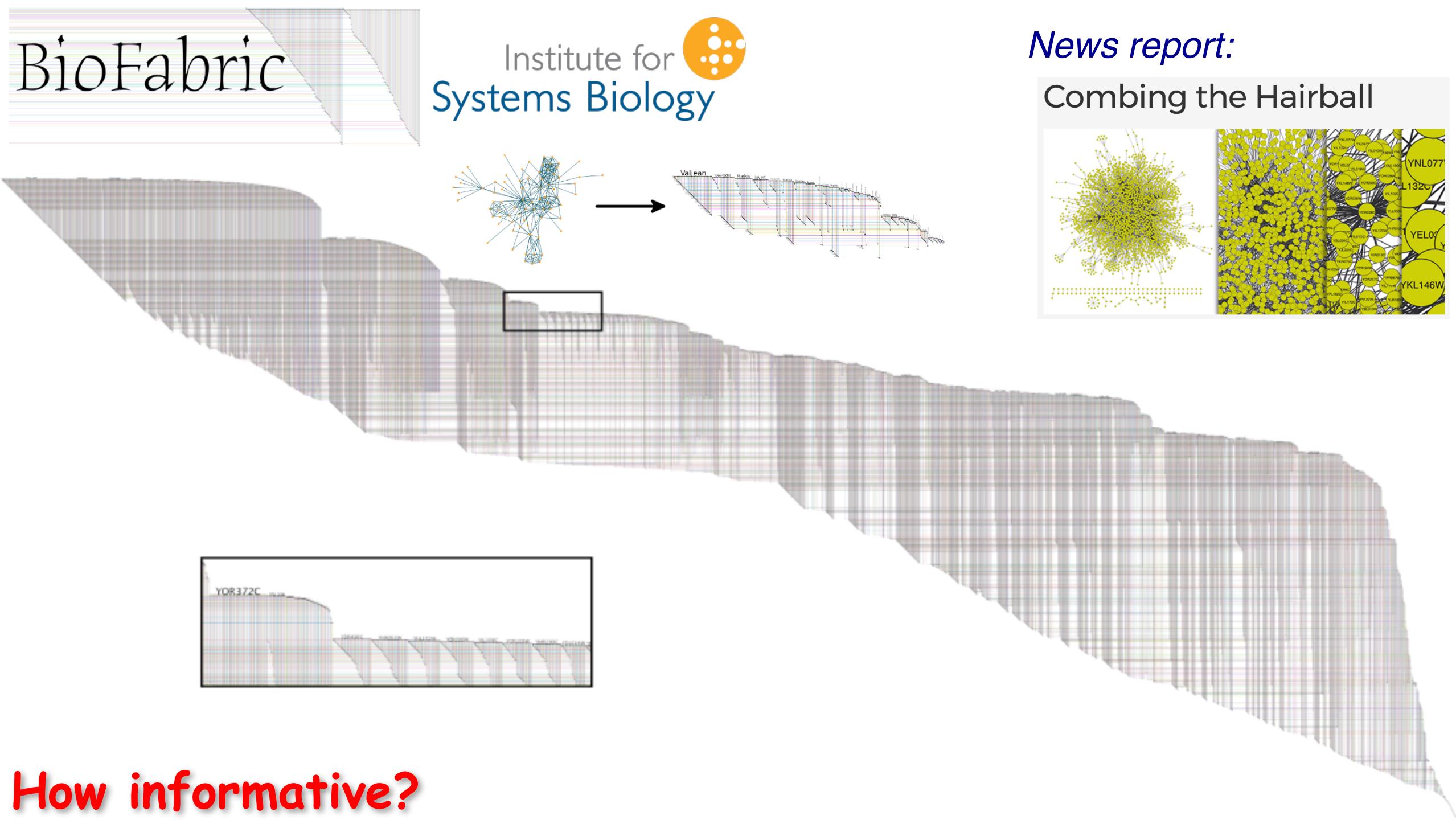
23 JANUARY 2004 VOL 303 SCIENCE



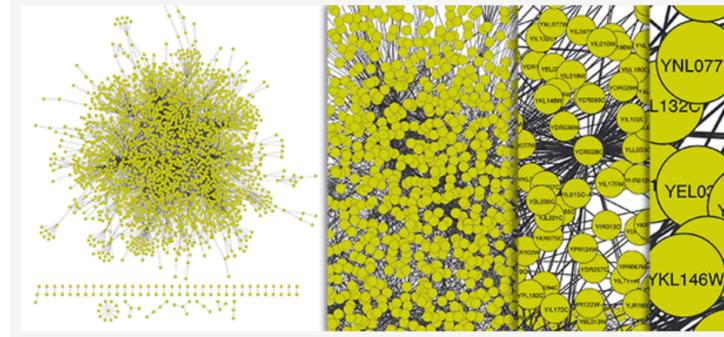
Interactome: the whole set
of molecular interactions in a
cell



How informative?



News report: Combing the Hairball



How informative?

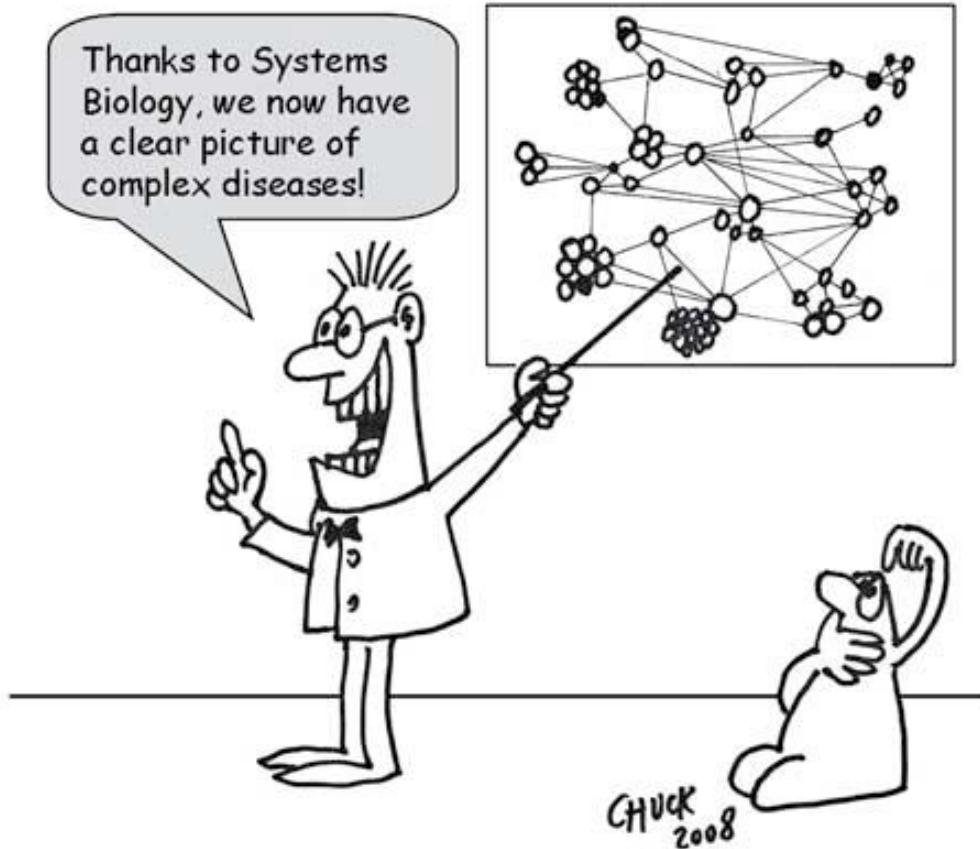
Pride and prejudice

Biologists:

- See quantitative science as just a tool
- Claim biological significance with weak statistical support
- Underestimate their lack of (S)TEM

Quantitative scientists:

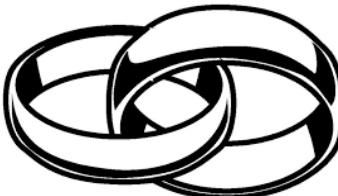
- Algorithm more important than the biological question
- Biological systems can be modeled by a set of simple equations (hard time dealing with biological variability)
- Find biology too descriptive and its conclusions far fetched



What makes good ~~computational biologists~~ biological data scientists

Drive to address
biological questions

(*i.e. strong interest in
biology*)



Fluency in
data analytics
(*e.g. algebra, statistics,
algorithms*)

Empathy!

*Genuine interest in the “other” side’s questions
and ability to understand how they think*

Our anti-black box efforts

- Make our ways of analysing data intelligible for (and therefore open for scrutiny by) collaborators
- Develop tools for empowering colleagues to perform similar analyses while understanding decisions needed to be made at their every stage



Thank you!

Questions?

BIOLOGY IS LARGELY SOLVED.
DNA IS THE SOURCE CODE
FOR OUR BODIES. NOW THAT
GENE SEQUENCING IS EASY,
WE JUST HAVE TO READ IT.

IT'S NOT JUST "SOURCE
CODE." THERE'S A TON
OF FEEDBACK AND
EXTERNAL PROCESSING.



BUT EVEN IF IT WERE, DNA IS THE
RESULT OF THE MOST AGGRESSIVE
OPTIMIZATION PROCESS IN THE
UNIVERSE, RUNNING IN PARALLEL
AT EVERY LEVEL, IN EVERY LIVING
THING, FOR FOUR BILLION YEARS.

IT'S STILL JUST CODE.



OK, TRY OPENING GOOGLE.COM
AND CLICKING "VIEW SOURCE."

OK, I... OH MY GOD.

THAT'S JUST A FEW YEARS OF
OPTIMIZATION BY GOOGLE DEVs.
DNA IS THOUSANDS OF TIMES
LONGER AND WAY, WAY WORSE.

WOW, BIOLOGY
IS IMPOSSIBLE.



*Genomic science is wonderful in that it brings together representatives of so many disciplines - clinicians, bench biologists, statisticians, bioinformatics scientists - all of whom tend to consider the others **intellectual peasants**.*

Zak Cohane

Prof. of Biomedical Informatics @ Harvard

