



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

Australian Institute for  
Bioengineering and Nanotechnology

# Gene Expression – Part 1

## Technologies for Transcriptomics

Associate Prof Jess Mar

Australian Institute for Bioengineering &  
Nanotechnology Level 4 West

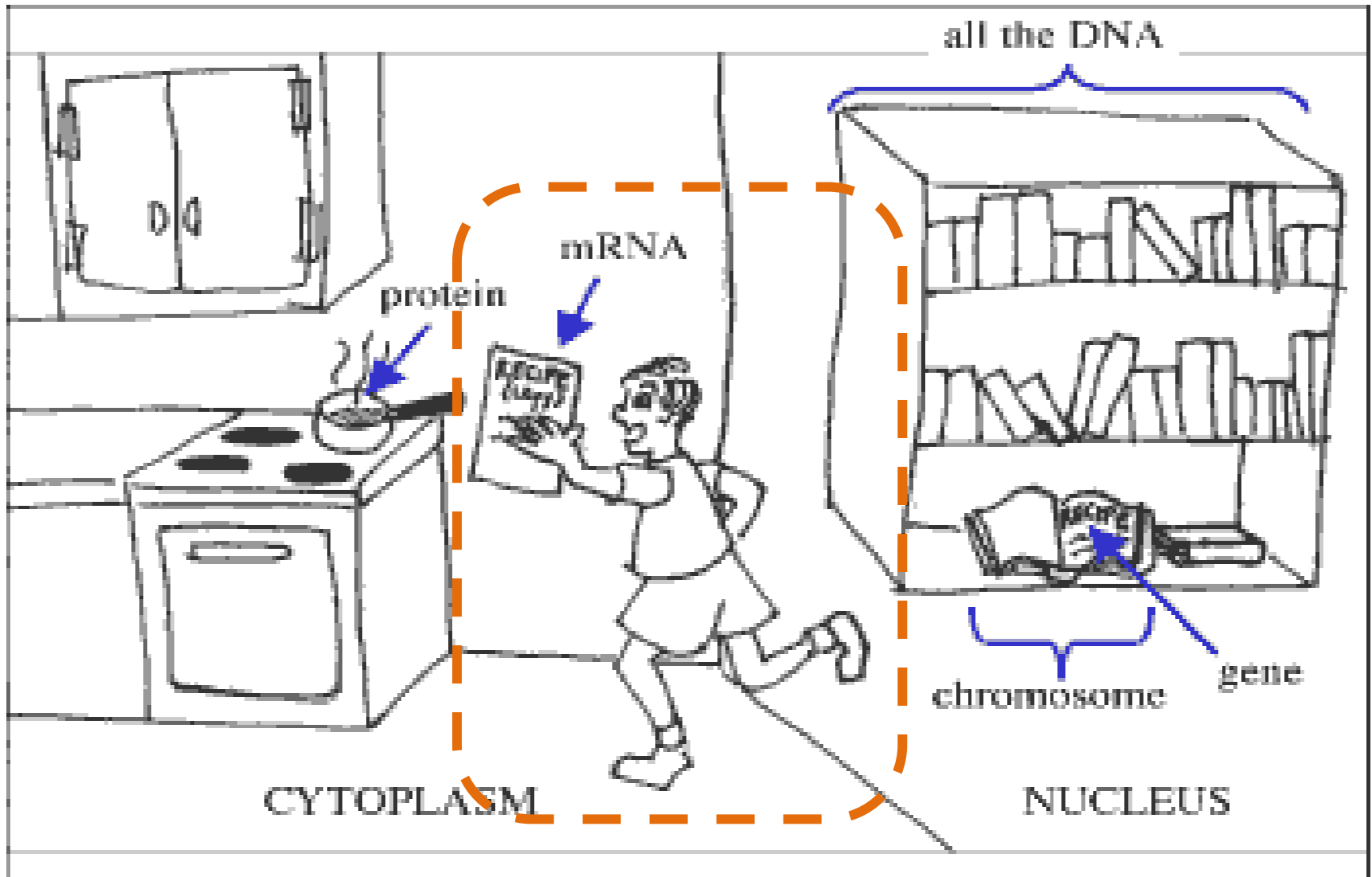
[jmar@uq.edu.au](mailto:jmar@uq.edu.au)

<https://aibn.uq.edu.au/mar>

 @jessicacmar

SCIE2100/BINF6000 – Semester 1, 2021

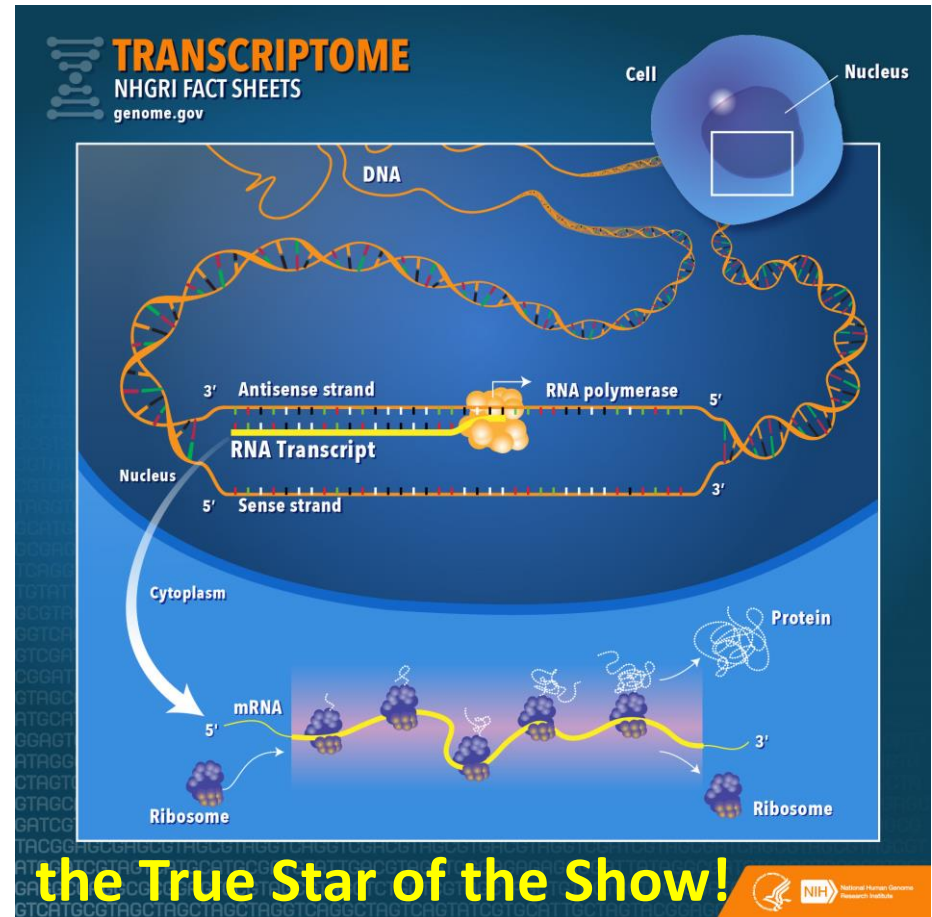
# ***The Central Dogma of Molecular Biology*** **has guided our interpretation of the genome**



# Introducing the Transcriptome

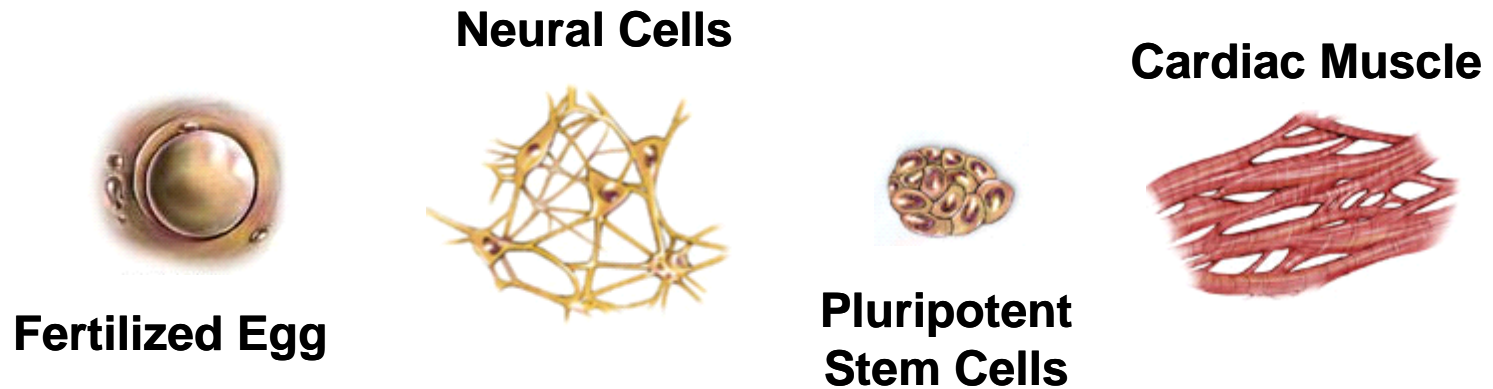
**Definition: the transcriptome is defined as the entire set of transcripts produced inside a cell which is subject to a developmental stage or physiological condition.**

- Transcripts represents RNA molecules including **mRNA, rRNA, tRNA, and other non-coding RNA** like lncRNA & microRNA.
- Genes are expressed or transcribed in response to a specific **cell stage** or **physiological condition**.
- Understanding transcriptional changes represents a window into figuring out the functional rules of the genome.



# The human body is made up of different kinds of cells

Within a single human, over 200 highly-specialized cell types exist!



Almost all cell types share the same genome.

**Epigenetic modification** and **transcription factor networks** generate the mechanism for cell type-specific diversity.

A cell type's unique properties are regulated by its transcriptional program.



# Understanding DNA variation is important but not the complete picture

*In the context of crime, cats, or anything else (equally) important!*

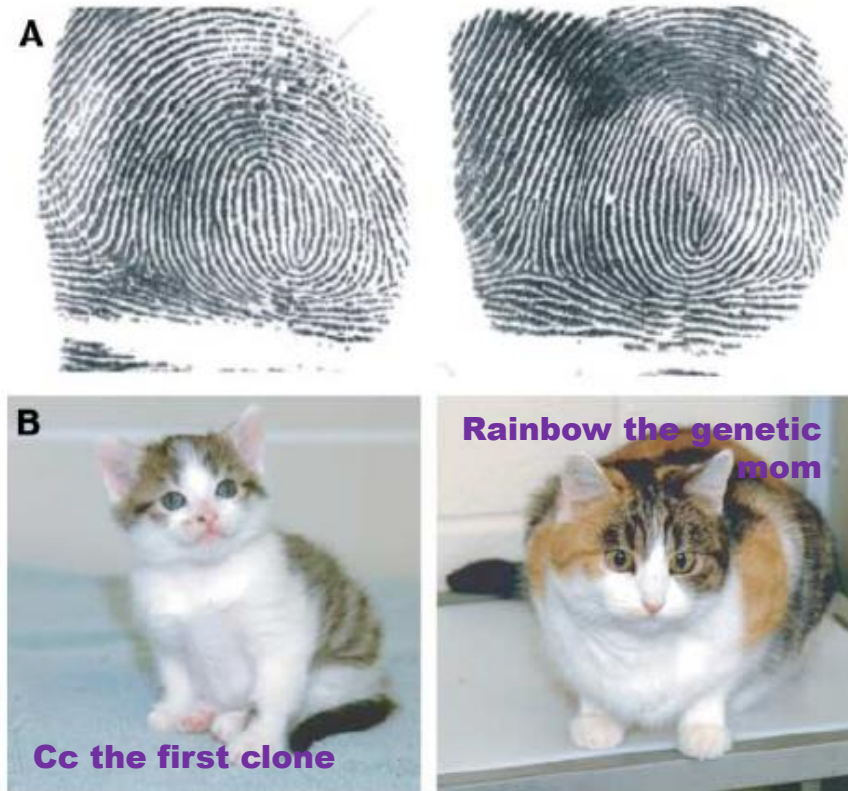
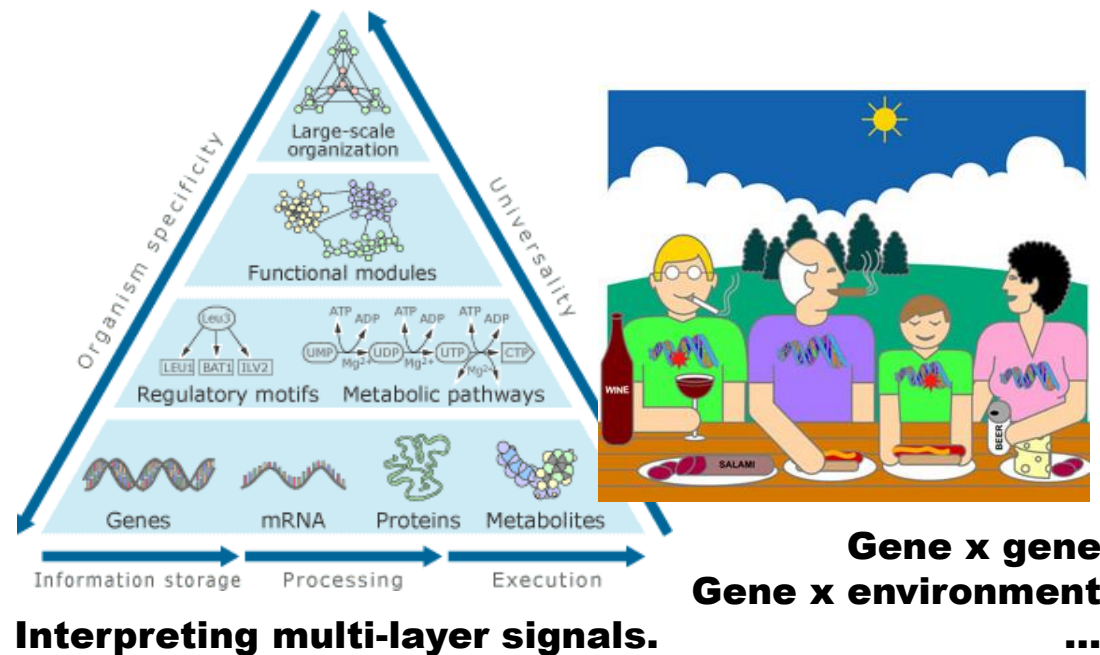


Fig. 1. Examples of possible stochastic influences on phenotype. (A) The fingerprints of identical twins are readily distinguished on close examination. Reprinted from (37) with permission from Elsevier. (B) Cc, the first cloned cat (left) and Rainbow, Cc's genetic mother (right), display different coat patterns and personalities (38). Photo credit, College of Veterinary Medicine and Biomedical Sciences, Texas A&M University.

# Biology is a complex system

The more we sequence, the more we recognize that there are **MANY** factors that contribute to how biology is regulated!

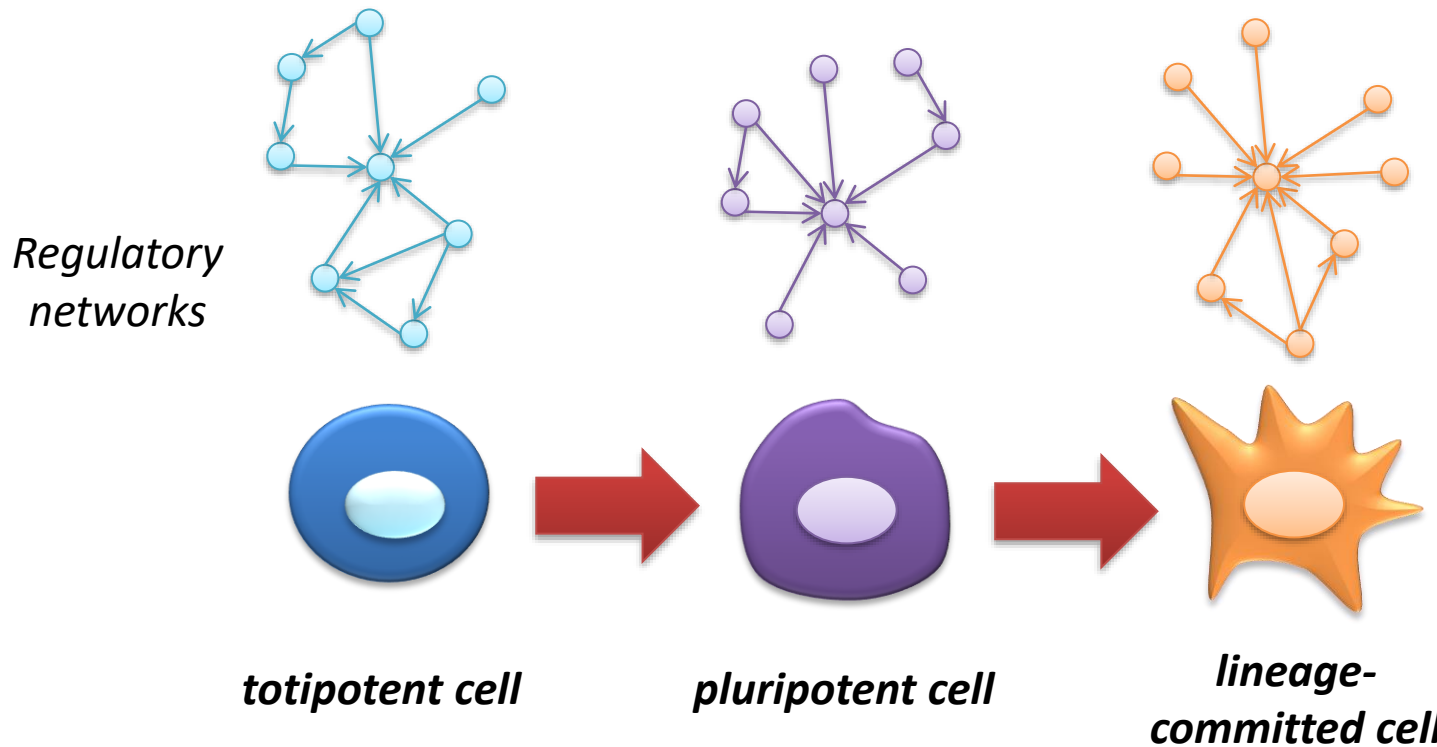


The goal of systems biology is to make sense of this complexity by developing predictive models and generating hypotheses that can be experimentally validated.

# Cellular phenotypes are controlled by gene regulatory networks

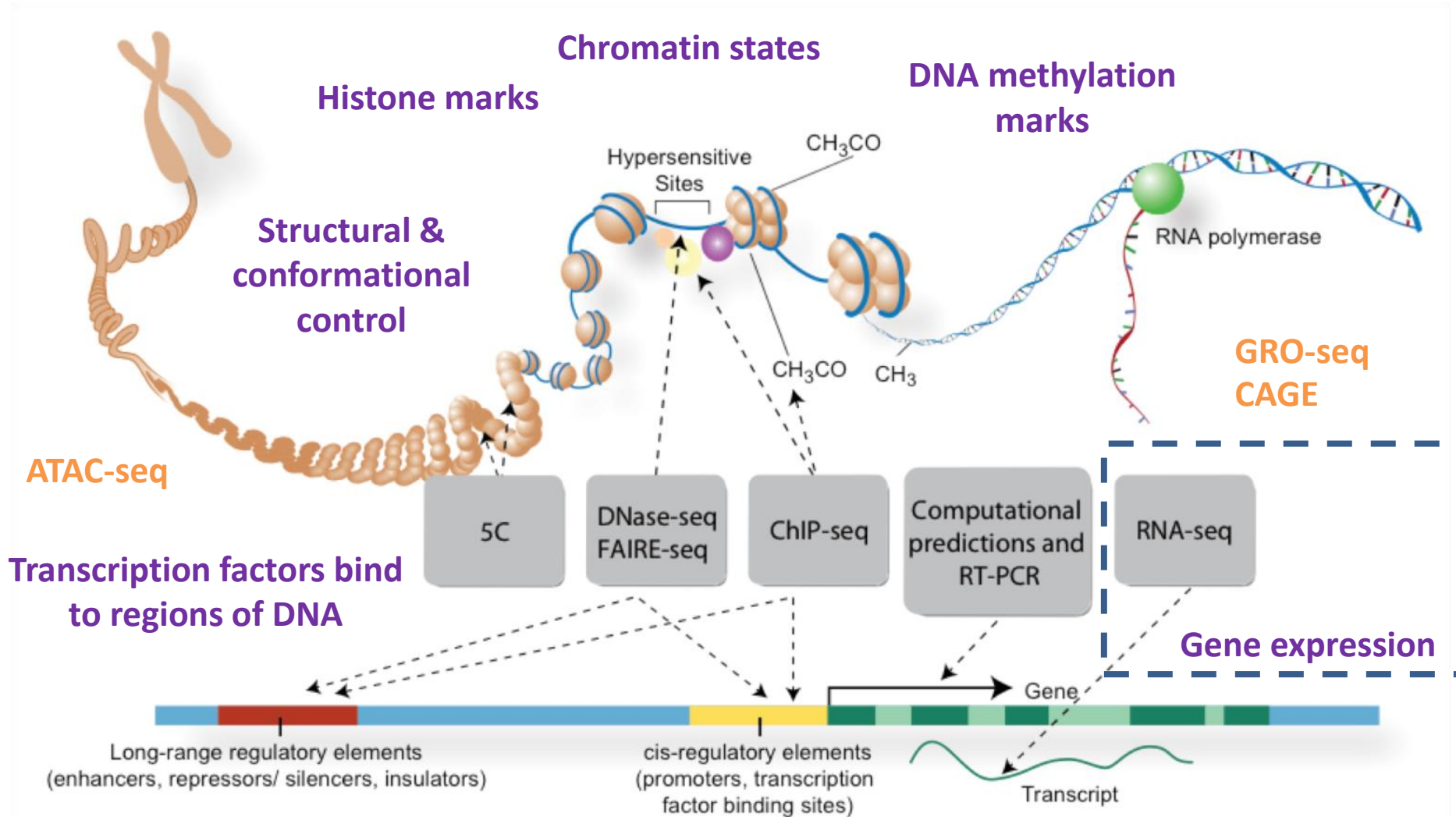
Genes work together in coordinated context-specific interactions that regulate a cellular phenotype.

Networks are graphical (graph + probability) models that capture our knowledge of which interactions define a phenotype.



*Our job as bioinformaticians is to integrate datasets to infer what these network models look like.*

# Improvements in technology continue to show that the genome is a very complex entity!

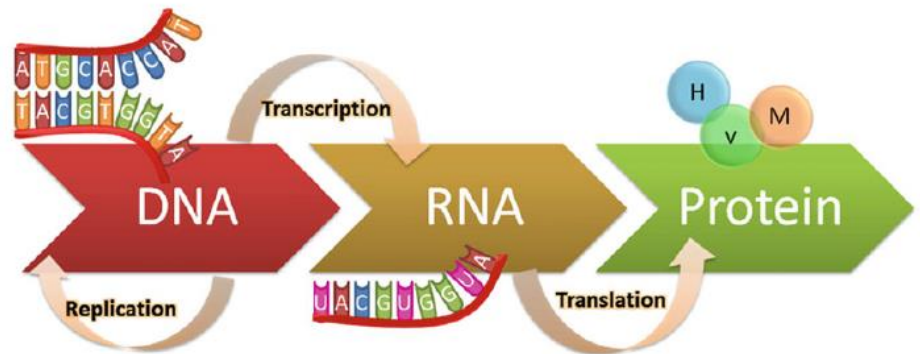




# The rationale for studying RNA

- RNA is another biological molecule, just like DNA.
- A protein-coding gene is made into a transcript (e.g. mRNA) before being turned into a protein inside a cell.
- Transcripts are copies of a gene, just like photocopies of a document.
- Cells order mRNA copies of a gene to be made when that gene is required for a specific task or purpose.
- The more copies (number of mRNAs), the more this gene is in demand!

**Studying mRNA counts is informative for understanding how biological processes are regulated.**



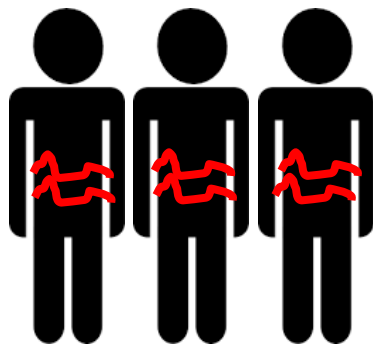
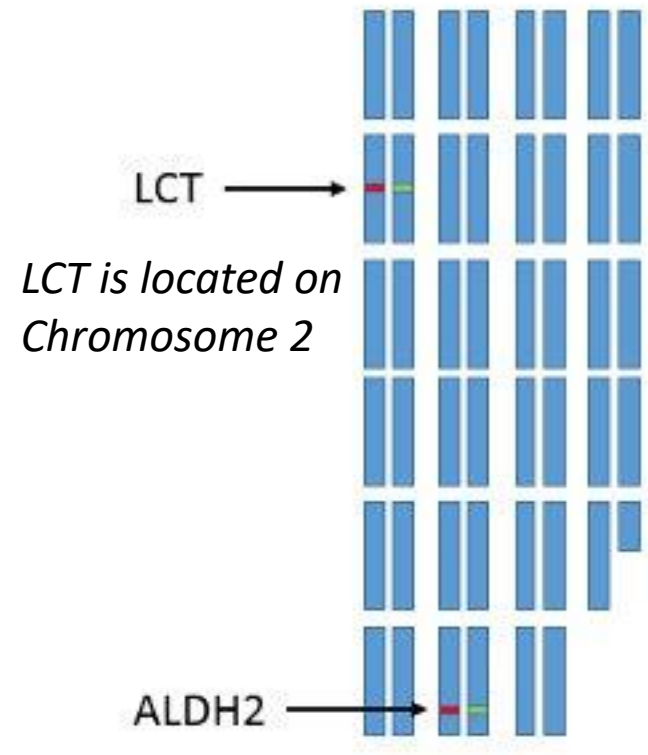


# Measuring DNA variants produces binary or categorical types of data

Example: Lactose intolerance is a condition where people are unable to digest lactose (a sugar found in milk and dairy products).

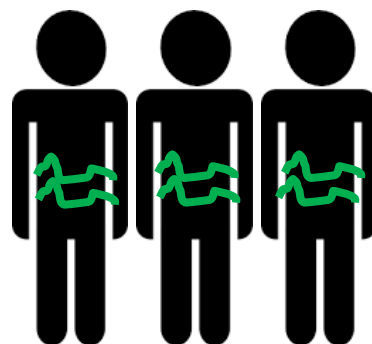
Lactose is broken up by an enzyme called lactase.

The LCT gene is responsible for making lactase.



*Homozygous  
Recessive  
(inherited lactose  
intolerance)*

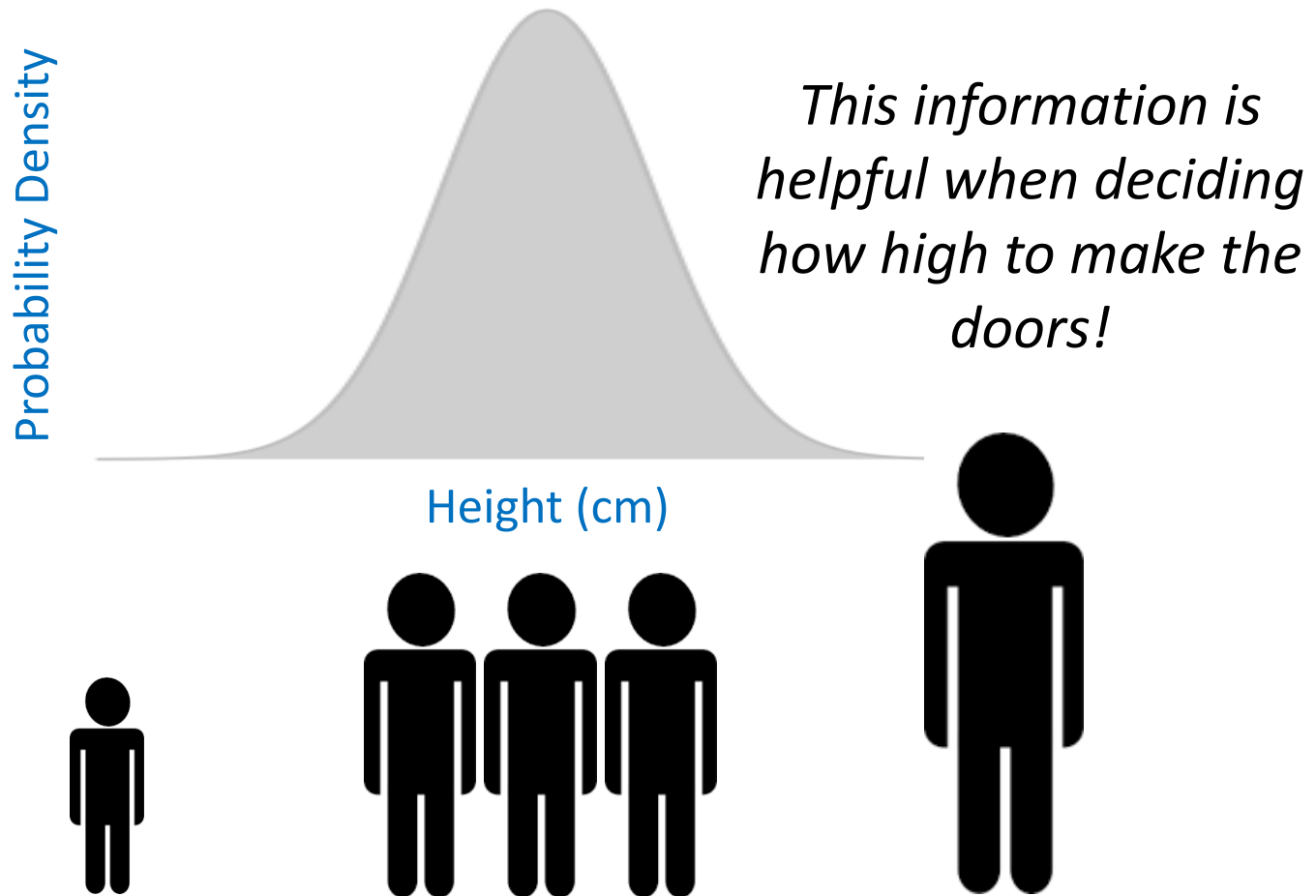
**vs**



*Heterozygous  
(usually no  
symptoms)*

# Measuring certain biometric variables produces continuous types of data

What is the average height of the student population at the University of Queensland?

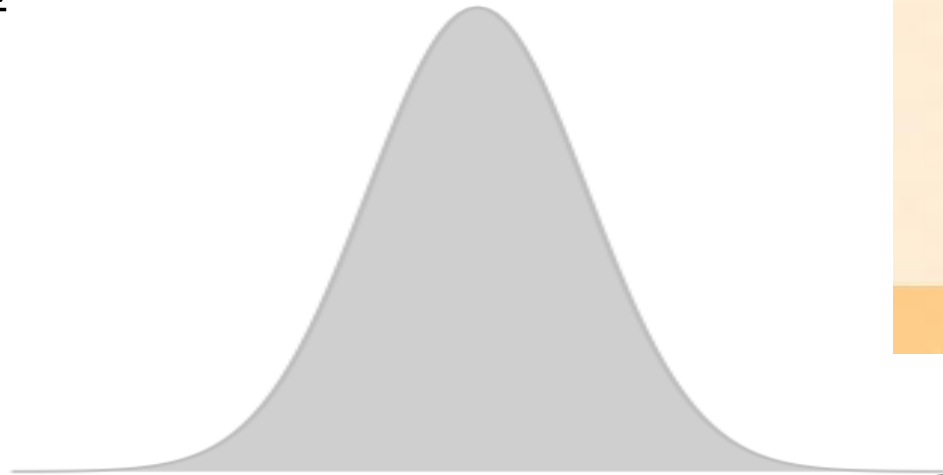


# Gene expression is a continuous (count) variable

Continuous data necessitates a different set of statistical analysis compared to binary/categorical data like DNA variation (see Part 2 next week!).

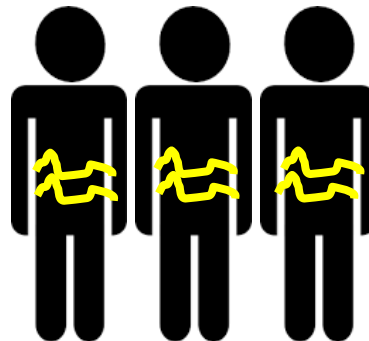
LCT gene example:

Probability Density



LCT gene expression

*Counting the number of  
LCT mRNA molecules.*



# Some technologies for measuring gene expression

- ***One to a few genes at a time***
  - Northern blotting
  - PCR-based approaches – e.g. quantitative RT-PCR
- ***Lots of genes in one go!***
  - Fluidigm (hundreds)
  - **DNA microarrays** (thousands to tens thousands)
  - Serial Analysis of Gene Expression (tens thousands)
  - Chromium 10X sequencing (thousands)
- ***All genes in one go!***
  - **RNA-sequencing**
  - Capped Analysis of Gene Expression (CAGE)



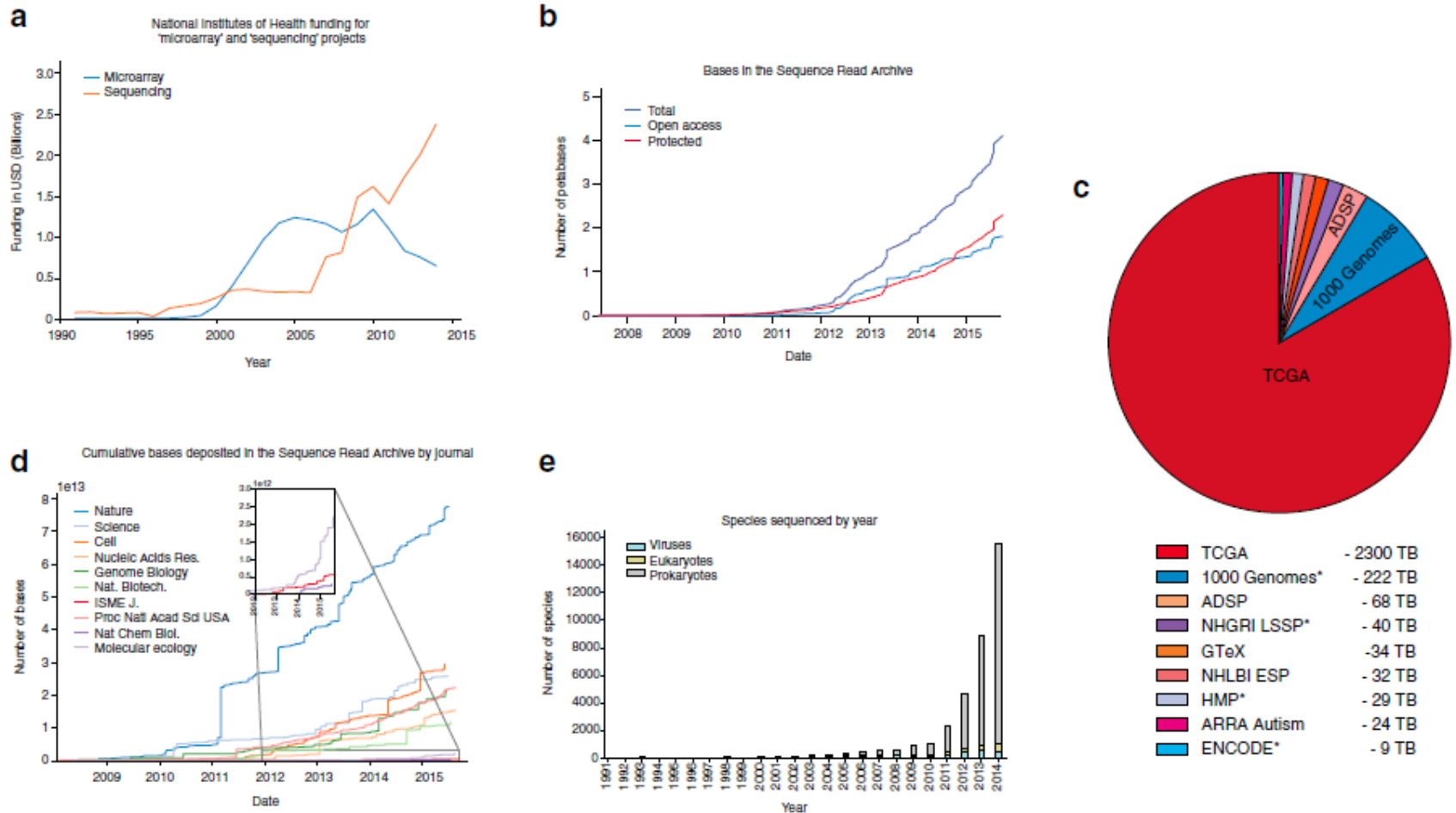


# Questions Addressed Today (& Next Week)

- What are the most common platforms for collecting high-throughput gene expression data?
- What are the key steps in analyzing RNA-sequencing and microarray data?
- How can we learn about biology through analyzing gene expression data?

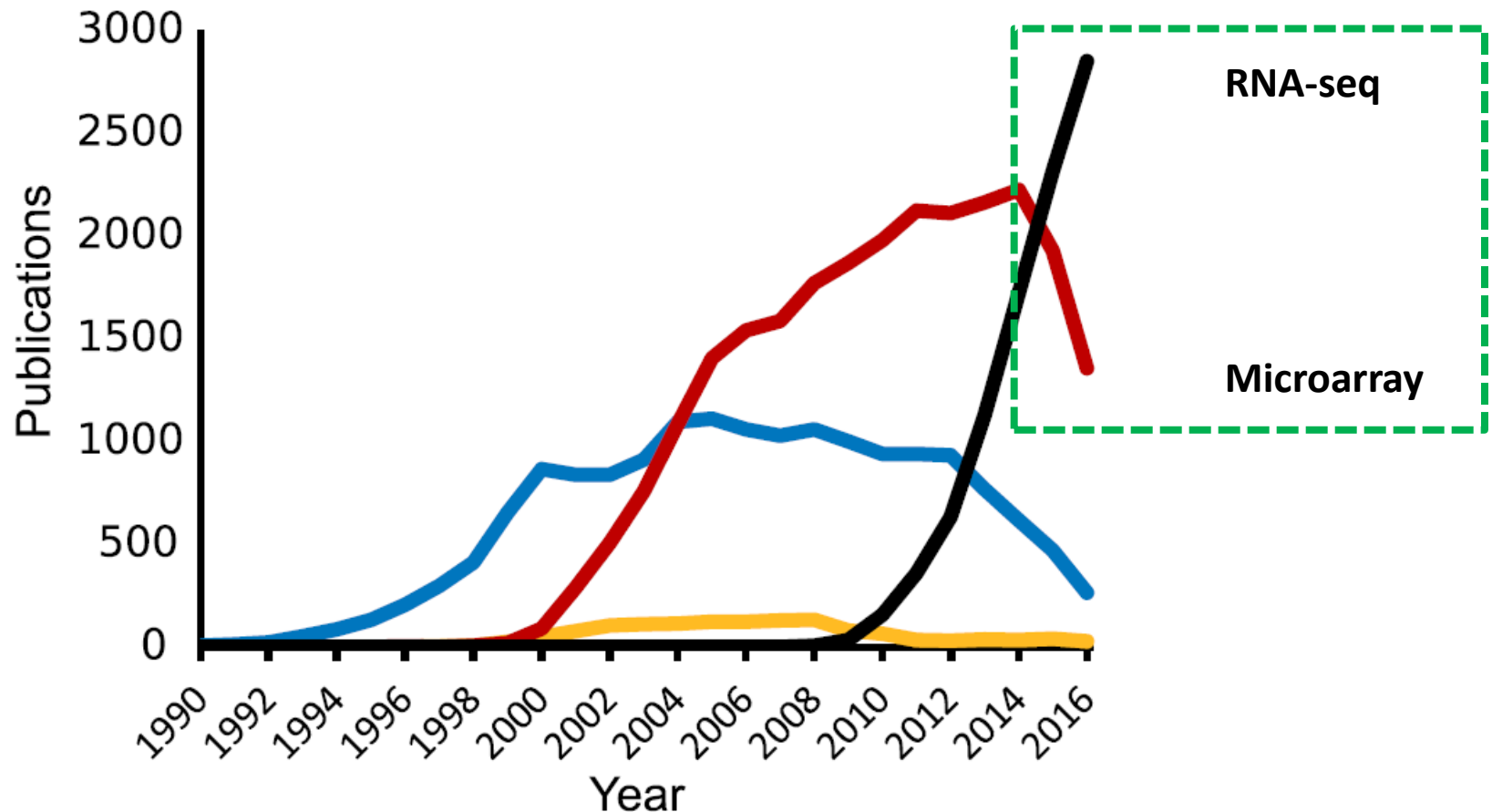
Let's start with technology

# Sequencing technologies are the driving force behind the big data revolution



Muir et al. (2016). *Genome Biology*. The real cost of sequencing: scaling computation to keep pace with data generation.

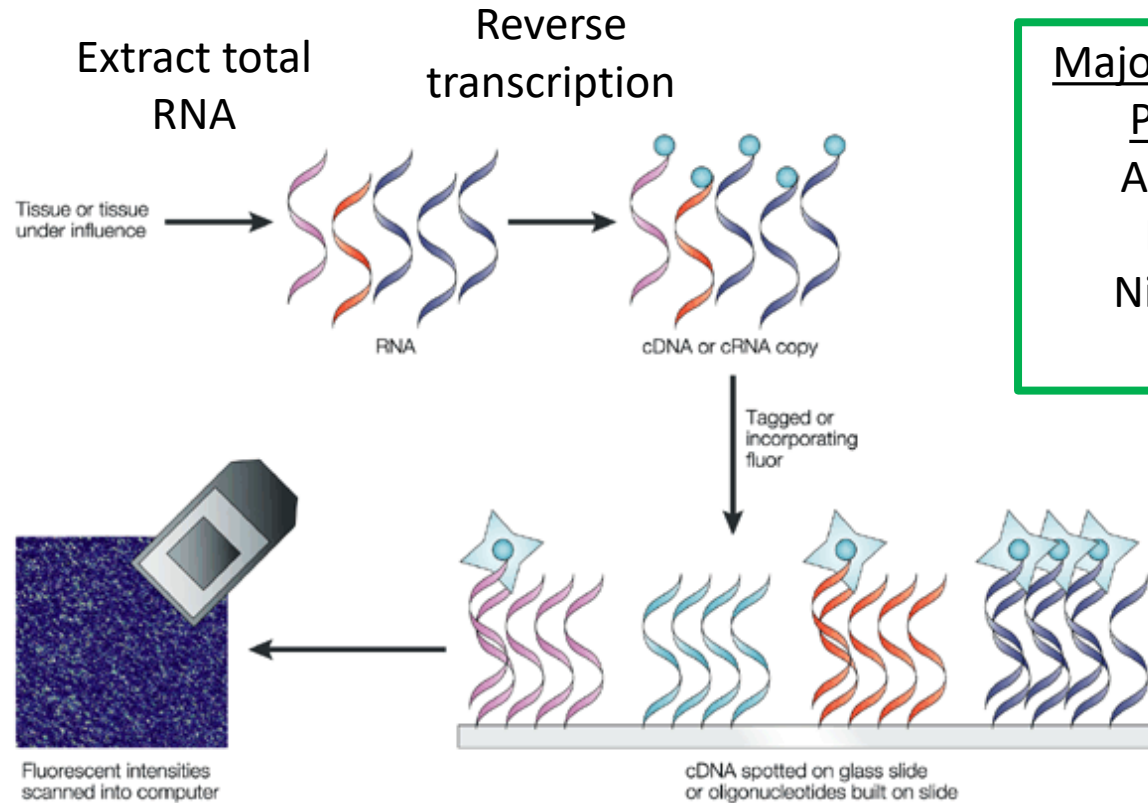
# RNA-sequencing is the heavy-weight amongst transcriptomics methods



**Fig 1. Transcriptomics method use over time.** Published papers since 1990, referring to RNA sequencing (black), RNA microarray (red), expressed sequence tag (blue), and serial/cap analysis of gene expression (yellow)[12].

<https://doi.org/10.1371/journal.pcbi.1005457.g001>

# But big data really started with microarrays...



## Major Microarray

### Platforms

Affymetrix

Illumina

NimbleGen

Agilent

Great (short) video intro!

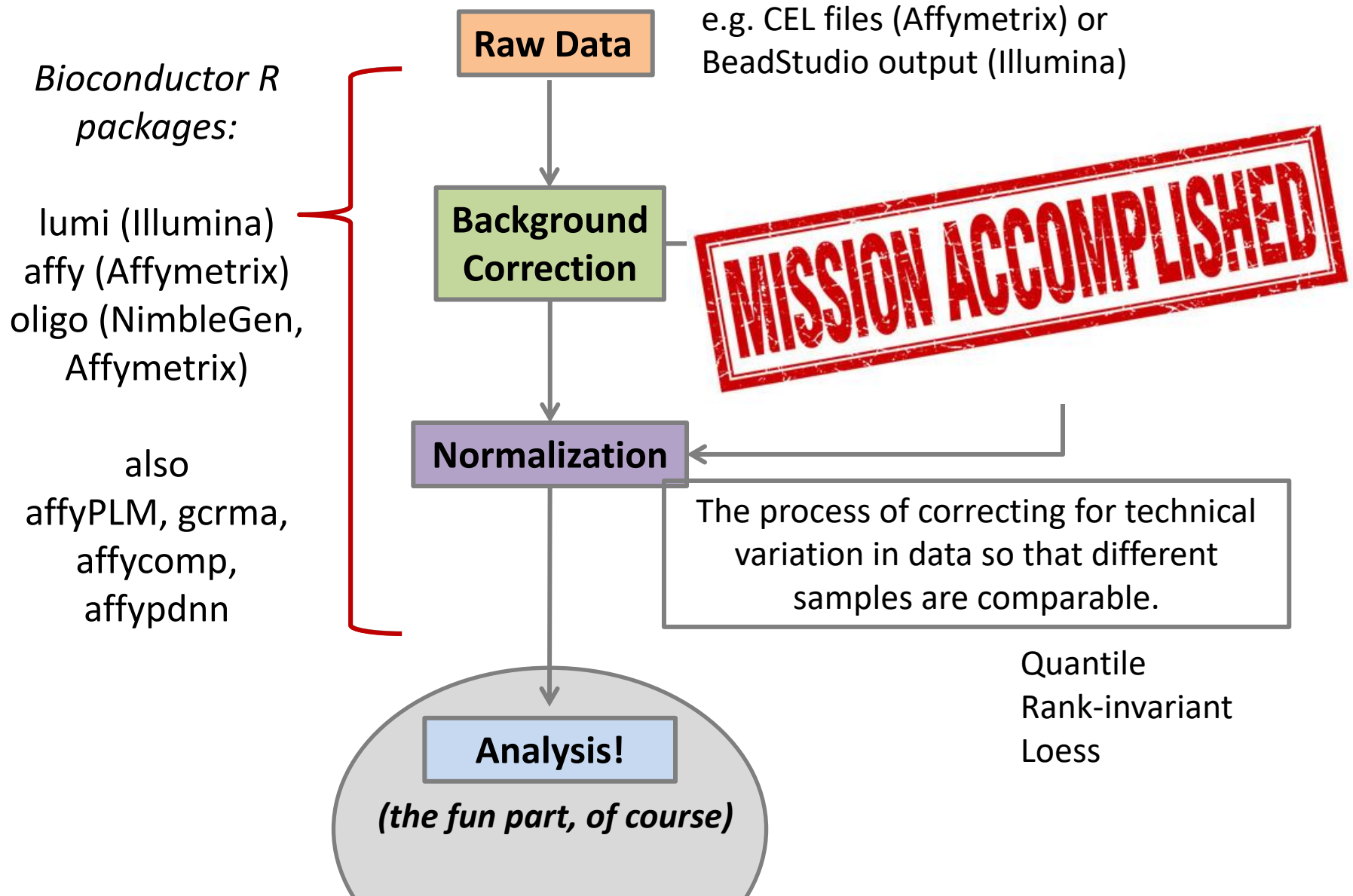
<https://www.youtube.com/watch?v=0ATUjAxNf6U>

Nature Reviews | Drug Discovery

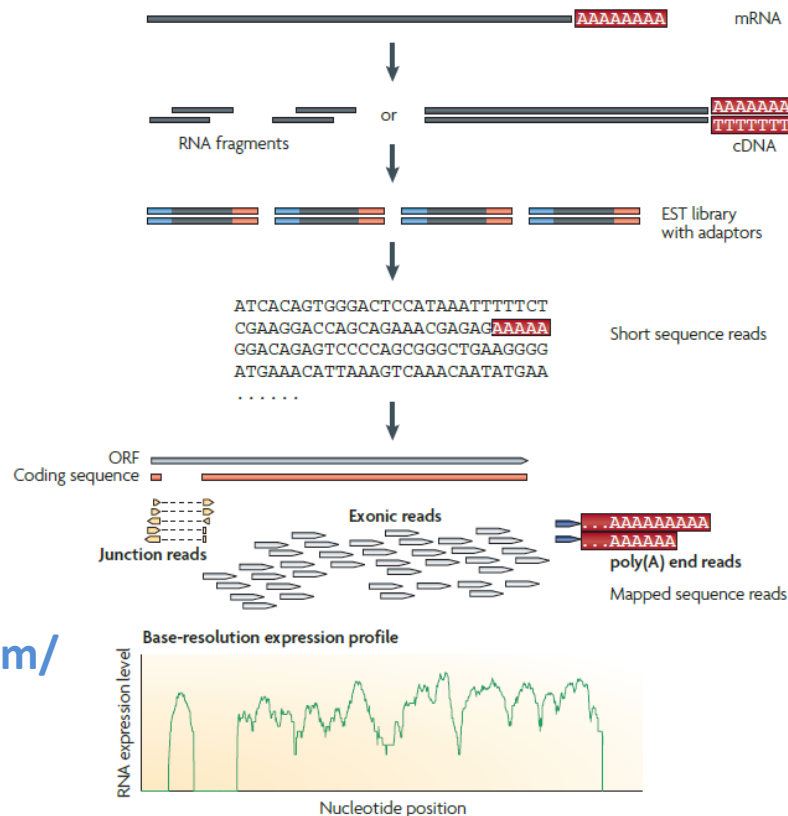
Butte. (2002). *Nat Rev Drug Discovery*. The use and analysis of microarray data.



# Standard pipeline for microarray pre-processing



# RNA-sequencing technology basics

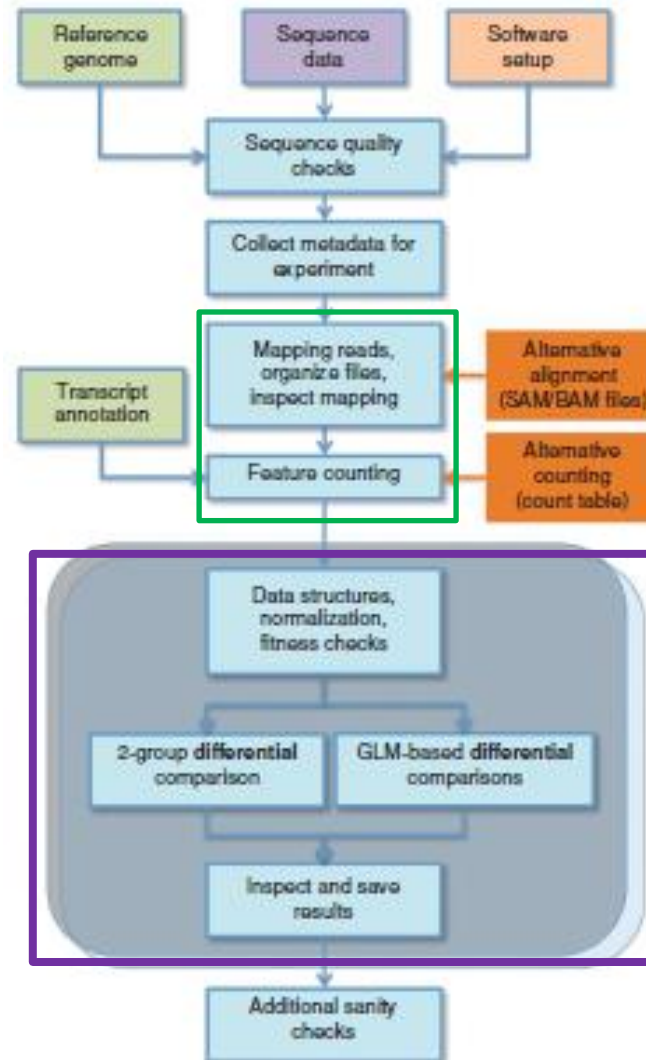


Illumina video:  
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Figure 1 | **A typical RNA-Seq experiment.** Briefly, long RNAs are first converted into a library of cDNA fragments through either RNA fragmentation or DNA fragmentation (see main text). Sequencing adaptors (blue) are subsequently added to each cDNA fragment and a short sequence is obtained from each cDNA using high-throughput sequencing technology. The resulting sequence reads are aligned with the reference genome or transcriptome, and classified as three types: exonic reads, junction reads and poly(A) end-reads. These three types are used to generate a base-resolution expression profile for each gene, as illustrated at the bottom; a yeast ORF with one intron is shown.

Wang et al. (2009). *Nature Reviews Genetics*. RNA-seq: a revolutionary tool for transcriptomics.

# RNA-seq analysis pipeline (simplified version)



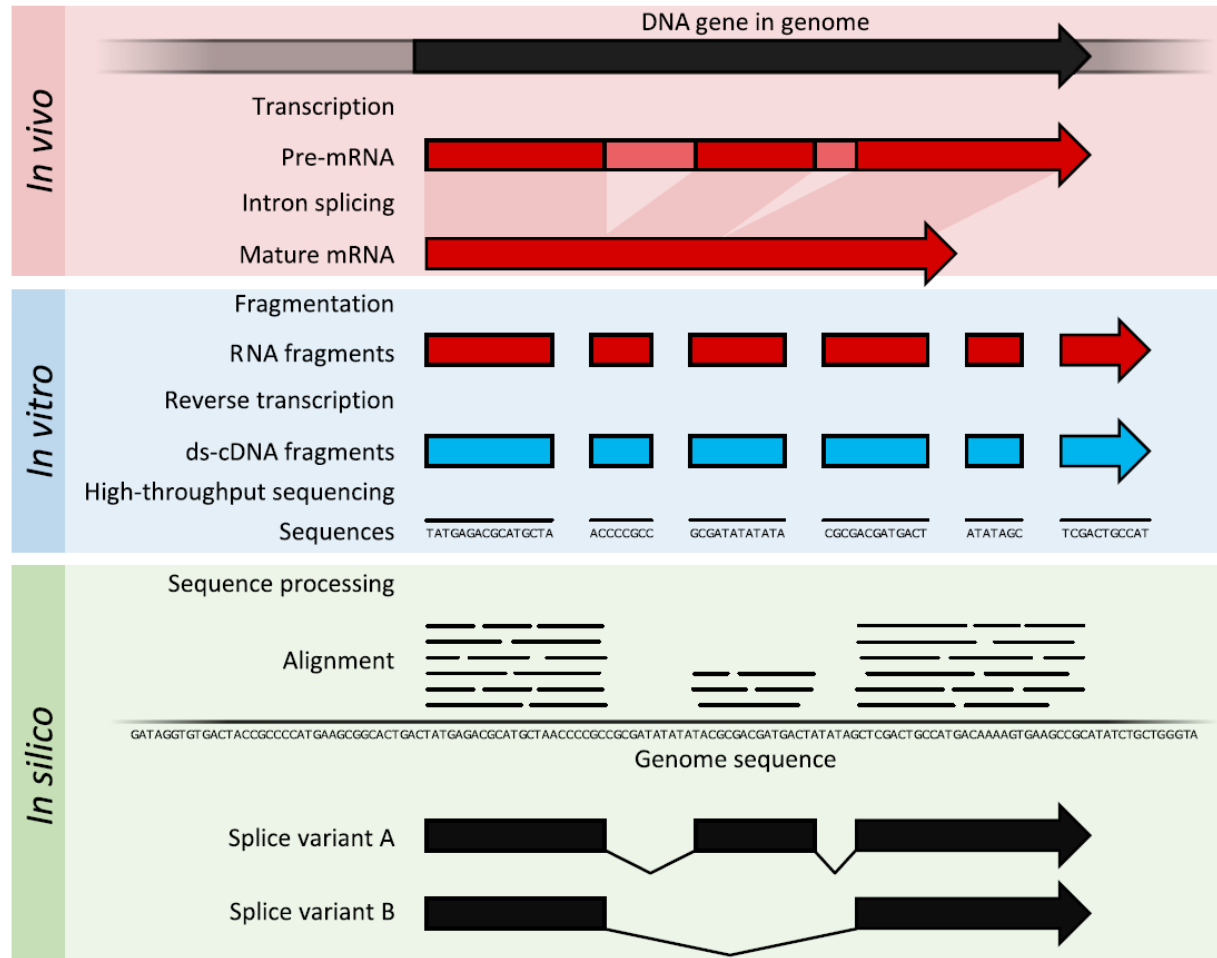
*Bioconductor R*  
*packages:*

edgeR

DESeq, DESeq2

Specialized Platforms:  
monocle (single cell)

# RNA-sequencing: breaking down the steps!



**Fig 4. Summary of RNA sequencing.** Within the organisms, genes are transcribed and spliced (in eukaryotes) to produce mature mRNA transcripts (red). The mRNA is extracted from the organism, fragmented and copied into stable double-stranded-cDNA (ds-cDNA; blue). The ds-cDNA is sequenced using [high-throughput](#), short-read sequencing methods. These sequences can then be [aligned](#) to a reference genome sequence to reconstruct which genome regions were being transcribed. These data can be used to annotate where expressed genes are, their relative expression levels, and any alternative splice variants.

<https://doi.org/10.1371/journal.pcbi.1005457.g004>

# Microarray vs RNA-sequencing



Table 1. Comparison of contemporary methods [23] [24] [10].

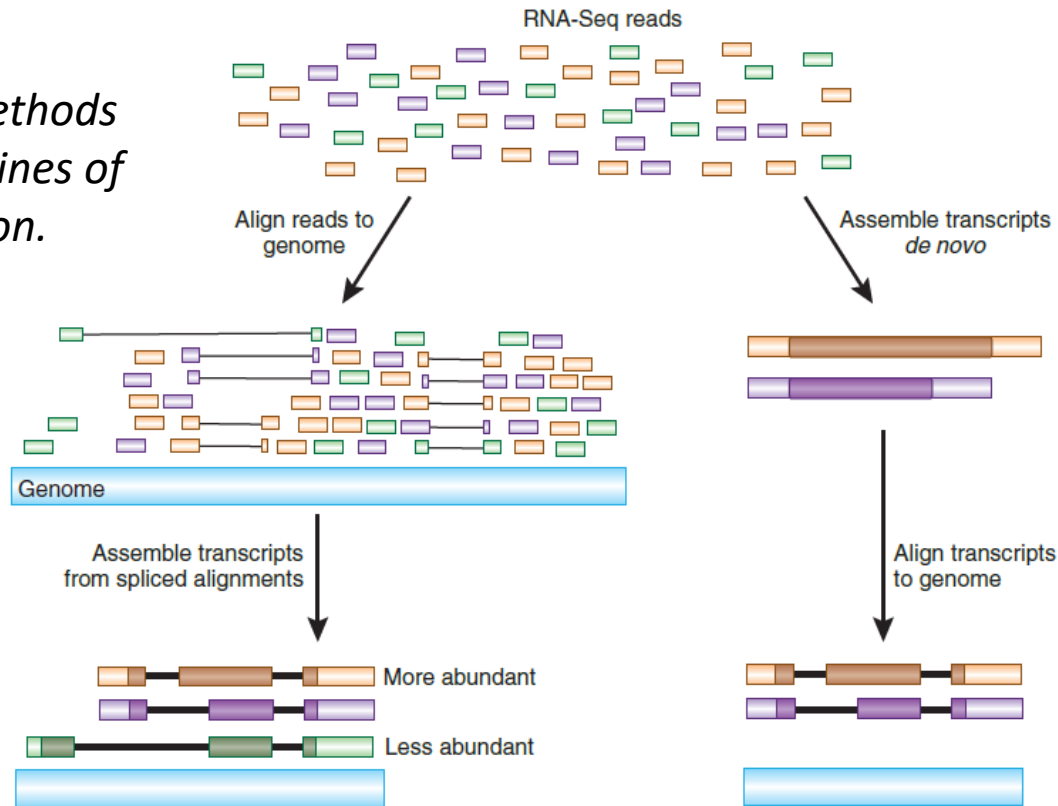
Method	RNA-Seq	Microarray
Throughput	High [10]	Higher [10]
Input RNA amount	Low ~ 1 <a href="#">ng</a> total RNA [25]	High ~ 1 $\mu$ g mRNA [26]
Labour intensity	High (sample preparation and data analysis) [10][23]	Low [10][23]
Prior knowledge	None required, though genome sequence useful [23]	Reference transcripts required for <a href="#">probes</a> [23]
<a href="#">Quantitation</a> accuracy	~90% (limited by sequence coverage) [27]	>90% (limited by fluorescence detection accuracy) [27]
Sequence resolution	Can detect <a href="#">SNPs</a> and splice variants (limited by sequencing accuracy of ~99%) [27]	Dedicated arrays can detect splice variants (limited by probe design and cross-hybridisation) [27]
Sensitivity	$10^{-6}$ (limited by sequence coverage) [27]	$10^{-3}$ (limited by fluorescence detection) [27]
Dynamic range	$>10^5$ (limited by sequence coverage) [28]	$10^3$ – $10^4$ (limited by fluorescence saturation) [28]
Technical reproducibility	>99% [29][30]	>99% [31][32]

- The size of the advantage depends on the **biological question** being asked!
- In general, key advantages (Table 1) have driven the uptick in RNA-seq versus microarray-based applications.



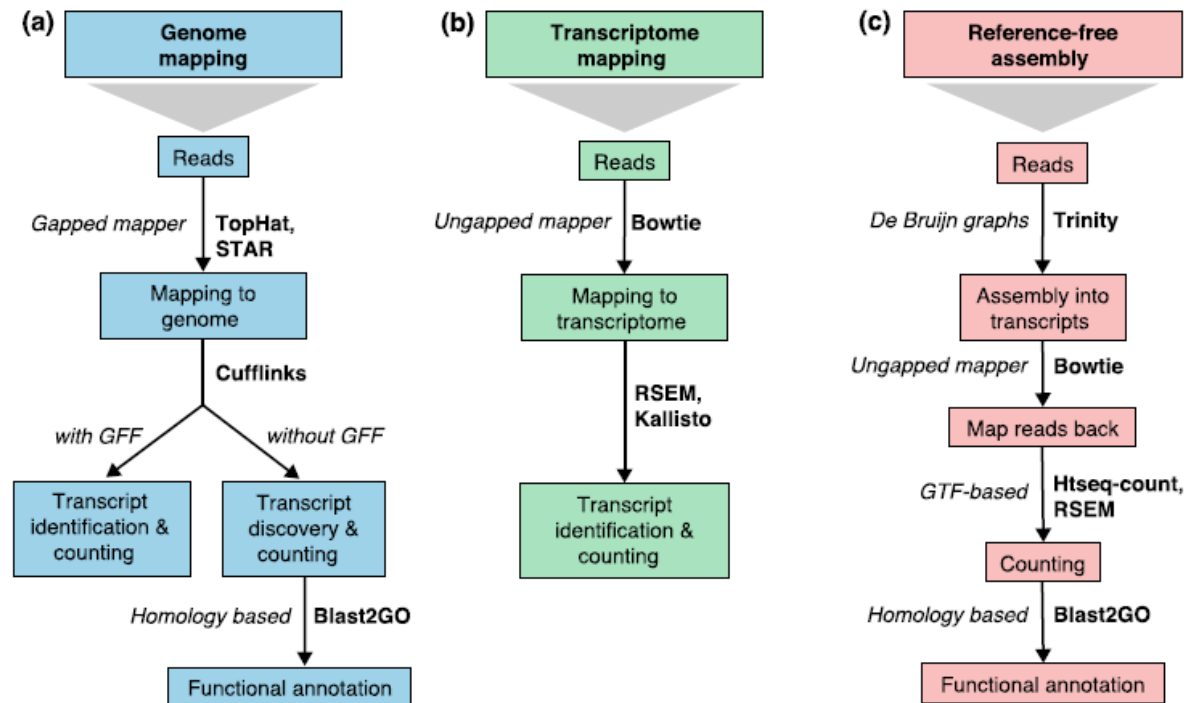
# Using RNA-seq to learn about the transcriptome

*Sequencing methods  
allow for new lines of  
investigation.*



Quantify differential expression, isoform abundance, non-coding RNAs, RNA editing, splicing, SNPs, etc.

# Despite the different applications, the preprocessing steps are generally similar



**Fig. 2** Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation file. Novel transcripts are then functionally annotated. **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis proceeds as in **(b)** followed by the functional annotation of the novel transcripts as in **(a)**. Representative software that can be used at each analysis step are indicated in *bold text*. Abbreviations: *GFF* General Feature Format, *GTF* gene transfer format, *RSEM* RNA-Seq by Expectation Maximization

# Technology platforms for next-generation sequencing

Table 2. Sequencing technology platforms commonly used for RNA-Seq [72][73].

Platform (Manufacturer)	Commercial release	Typical read length	Maximum throughput per run	Single read accuracy	RNA-Seq runs deposited in the NCBI SRA (Oct 2016) [74]
<a href="#">454</a> (Roche, Basel, Switzerland)	2005	700 bp	0.7 Gbp	99.9%	3548
<a href="#">Illumina</a> (Illumina, San Diego, CA, USA)	2006	50–300 bp	900 Gbp	99.9%	362903
<a href="#">SOLiD</a> (Thermo Fisher Scientific, Waltham, MA, USA)	2008	50 bp	320 Gbp	99.9%	7032
<a href="#">Ion Torrent</a> (Thermo Fisher Scientific, Waltham, MA, USA)	2010	400 bp	30 Gbp	98%	1953
<a href="#">PacBio</a> (Pacbio, Menlo Park, CA, USA)	2011	10,000 bp	2 Gbp	87%	160

NCBI, National Center for Biotechnology Information; SRA, Sequence Read Archive; RNA-Seq, RNA sequencing.

<https://doi.org/10.1371/journal.pcbi.1005457.t002>

New platforms and technologies continue to be released at rates that no publication can capture accurately.

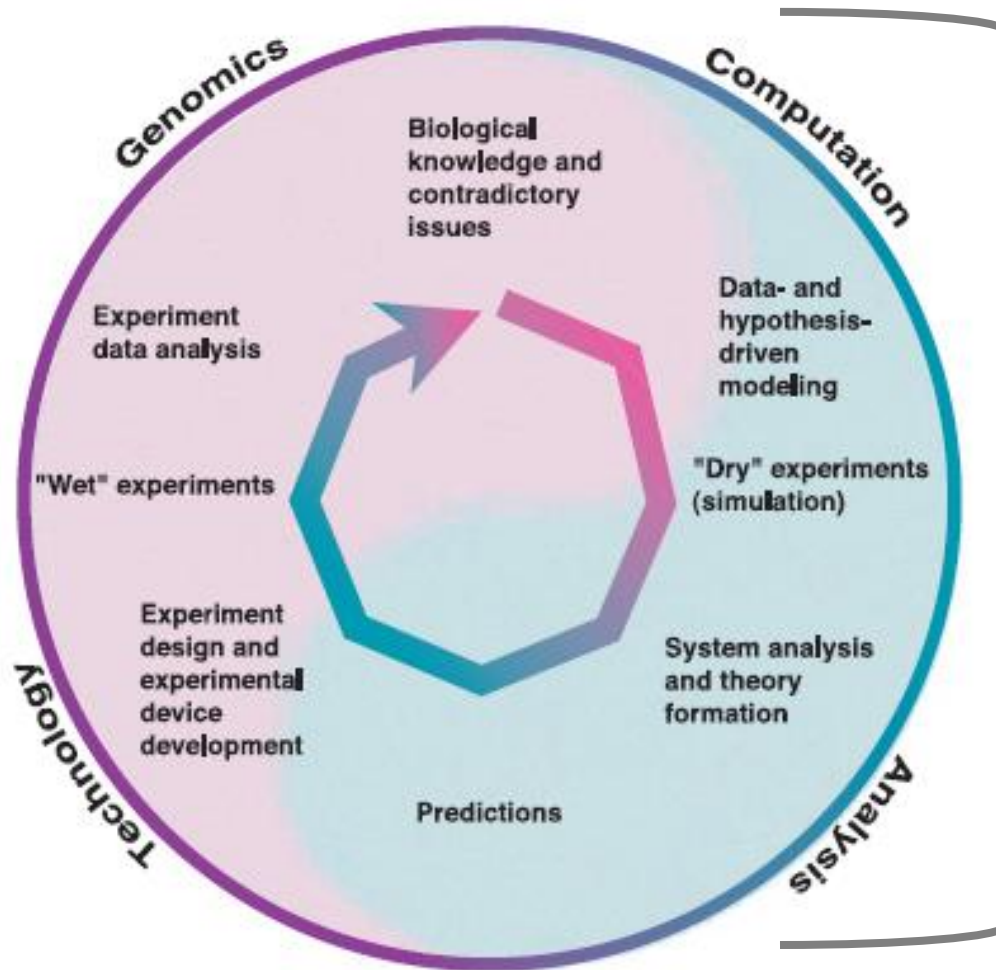
Lowe et al. (2017). *PLoS Computational Biology*. Transcriptomics technologies.

Resources to study gene expression

—

what's in this for you?

# How systems biology fits into the general scientific method



1. Identify a testable hypothesis.
2. Identify the data set(s) or experimental design under which this hypothesis can be tested.
3. Conduct statistical analysis.
4. Interpret results and acknowledge limitations.



# Got questions? Ask the data!

A wealth of gene expression data is publicly available in the form of repositories, the largest being GEO, and is hosted by the NCBI at the NIH.

## 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.



- All expression data stored in GEO follows a standard format called MIAME which was created to ensure the experimental details were available.
- Easy links to the paper that describes the study that generated the data.
- Data can be downloaded in normalized (good to go!), or raw formats.

**NCBI has many databases that store specific genomic information.**

dbSNP

dbGAP: database of Genotypes and Phenotypes

SRA (Short Read Archive)

<http://www.ncbi.nlm.nih.gov/>



Normalized metadata for the [Sequence Read Archive](#)

# Specialized databases continue to be added...

Find human RNA-seq samples <sup>?</sup> RESET

matching **all** of these terms: <sup>?</sup>

but **none** of these terms: <sup>?</sup>

Sample type: All cell line tissue primary cells stem cells in vitro differentiated cells iPS cell line

## Examples

- Find healthy liver tissue: require liver, exclude disease and treatment. Sample type: tissue.
- Find healthy, primary T-cells: require T cell, exclude disease and treatment. Sample type: primary cells.
- Find glioblastoma samples: require glioblastoma multiforme and brain.

Key: ○ Anatomy ○ Disease ○ Cell Line ○ Cell Type ○ Experimental Factor



<http://metasra.biostat.wisc.edu/>



# Social Genomics – Loneliness, Happiness and Science?!

An emerging area of social science deals with the intersection of happiness/loneliness and the impact on human health. More recently, this field has taken a quantitative molecular approach, giving rise to “social genomics”.

## Loneliness Is Bad For You, And This Study May Help Explain Why

Feeling lonely may trigger changes in our cells that could make us more susceptible to illness.

11/28/2015 08:53 am ET



Jacqueline Howard  
Senior Science Editor, The Huffington Post



EVGENIASH VIA GETTY IMAGES

Loneliness can affect the production of white blood cells in our bodies, study shows.

Forbes / Pharma & Healthcare

VE. INFORMATE.  
TOMA CONTROL.



NOV 24, 2015 @ 08:00 AM 15,913 VIEWS

## Loneliness Destroys Physical Health From The Inside Out



David DiSalvo  
CONTRIBUTOR

I write about science, technology and the cultural ripples of both.



FULL BIO >

Opinions expressed by Forbes Contributors are their own.

Loneliness can increase the risk of premature death in older adults by 8%, a new study supported by the National Institute of Aging shows.

What the research team found is that loneliness affects the expression of genes related to two critical physiological systems and increased cellular inflammation. CTRA simultaneously increasing the genetic expression of genes at the cellular level rather than the system level.

The longer someone experiences loneliness, the more genes related to white blood cells (aka inflammation) and inflammation. CTRA simultaneously increasing the genetic expression of genes at the cellular level rather than the system level.

The combination of the two effects is with a slow erosion of cellular health problems, some of which worsen over time.

The study also found that CTRA and CTRA gene expression more than a year later. In other words,

## The Physical Effects Of Loneliness Include Weakened Immune Systems, Premature Death

AFP/Relaxnews

Posted: 11/24/2015 10:50 am EST | Updated: 11/24/2015 10:59 am EST



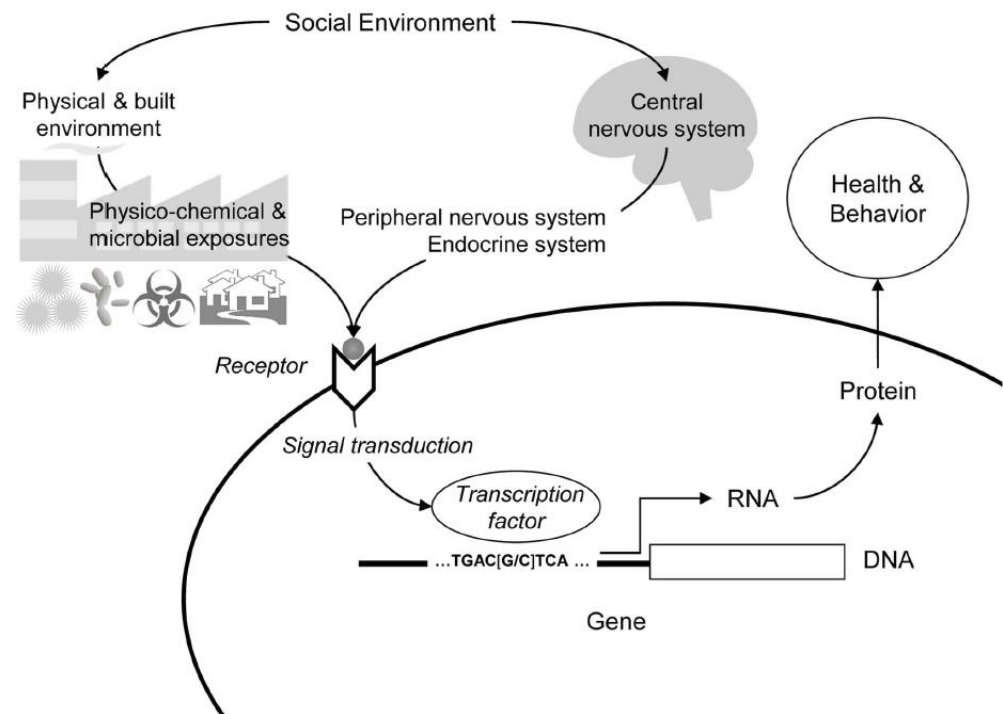
# Do Social-Environmental Factors Influence Our Gene Expression?

Proponents of human social genomics argue that social-environmental conditions can affect differential expression of genes in leukocytes.

Circulating leukocytes are a mixed population of cells, responsible for pathogen recognition, immune response and tissue repair.

Adverse social conditions (loneliness, bereavement, depression) have shown to affect leukocyte gene expression.

*How do social environments regulate immune function?*



**Figure 1. Social regulation of human gene expression.** Social environments can influence human gene expression via physicochemical processes (e.g., toxins, pollutants, and microbes) and psychological processes (e.g., experiences of threat or uncertainty) that trigger neural and endocrine responses (e.g., activation of the sympathetic nervous system). In both cases, biochemical mediators engage cellular receptor systems, which activate intracellular signal transduction pathways culminating in the activation (or repression) of transcription factors that proximally regulate the transcription of genes bearing response elements for that particular factor. The gene regulatory “wiring diagram” that maps specific biochemical signals to specific gene expression responses represents an evolved genomic program that was presumably adaptive under ancestral conditions but may have distinct maladaptive effects in the very different social environments of contemporary human life.  
doi:10.1371/journal.pgen.1004601.g001

**Cole SW. (2014).** Human Social Genomics. *PLoS Genetics*.

# Retrieving data from GEO is straightforward

Research

Open Access

## Social regulation of gene expression in human leukocytes

Steve W Cole<sup>\*†‡</sup>, Louise C Hawkey<sup>§</sup>, Jesusa M Arevalo<sup>\*</sup>, Caroline Y Sung<sup>†</sup>, Robert M Rose<sup>¶</sup> and John T Cacioppo<sup>§</sup>

Published: 13 September 2007

Genome **Biology** 2007, **8**:R189 (doi:10.1186/gb-2007-8-9-r189)

Received: 2 March 2007

Revised: 30 July 2007

Accepted: 13 September 2007

from the [Methods](#) section

### Social isolation

Subjectively experienced social isolation was assessed by the UCLA-R Loneliness scale [53] at each yearly study visit. Biological samples from 10 individuals who consistently scored in the top 15% of the loneliness distribution during study years 1, 2, and 3, and 10 individuals who consistently scored in the bottom 15% during years 1, 2, and 3, were selected for analysis after matching for age, gender, and ethnicity. Two samples from low-lonely individuals and four samples from high-lonely individuals yielded insufficient RNA for reliable gene expression assay, and analyses are thus based on fourteen individuals (eight low-lonely, six high-lonely). Objective social isolation was assessed by the social network index [54].

puncture sample and isolated by ficoll density gradient centrifugation; RNAlater/RNeasy, Qiagen, Valencia, CA, USA), and 5 µg of the resulting RNA was assayed using Affymetrix U133A high-density oligonucleotide arrays [58] in the UCLA DNA Microarray Core as previously described [41,59]. Robust multiarray averaging [60] was applied to quantify expression of the 22,283 assayed transcripts, and differentially expressed genes were identified as those showing ≥30% difference in mean expression levels in samples from high- versus low-lonely individuals (corresponding to a FDR of 10%) [59]. Functional characteristics of individual genes were identified through GO annotations, Gene References into Function annotations, and PubMed literature links retrieved through NCBI Entrez-Gene [61]. Functional commonalities among multiple differentially expressed genes were identified using GStat [62] with default stringency parameters (Benjamini FDR <0.10) [44]. A full list of differentially expressed genes is provided in Additional data file 1, and raw data are deposited in the National Center for Biotechnology Information Gene Expression Omnibus (GEO Series [GSE7148](#)).

To install any Bioconductor package:

```
> source("http://www.bioconductor.org/biocLite.R")
> biocLite("GEOquery")
```



# The Bioconductor Project – A Bioinformatics Standard

- This project has become the standard repository for R software that deals with all things **bio**.
- A big theme of Bioconductor has been the standardization of data classes to make analysis of –omic data easier, more robust and **more reproducible**.
- The project makes available packages that deal with:
  - Annotation
  - Statistical Methods
  - Pre-processing Approaches
- *Vignettes will change your life!*

<http://bioconductor.org>



## 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, [1024 software packages](#), and an active user community. Bioconductor is also available as an [AMI](#) (Amazon Machine Image) and a series of [Docker](#) images.

## News

- [Bioconductor 3.1 is released](#)
- *Nature Methods* Orchestrating high-throughput genomic analysis with Bioconductor ([abstract](#); full-text free with registration) and other recent [literature citations](#).
- Read our latest [newsletter](#).
- Updated [course material](#) and [videos](#).
- Use the [support site](#) to get help installing, learning and using Bioconductor.

# Lecture Summary

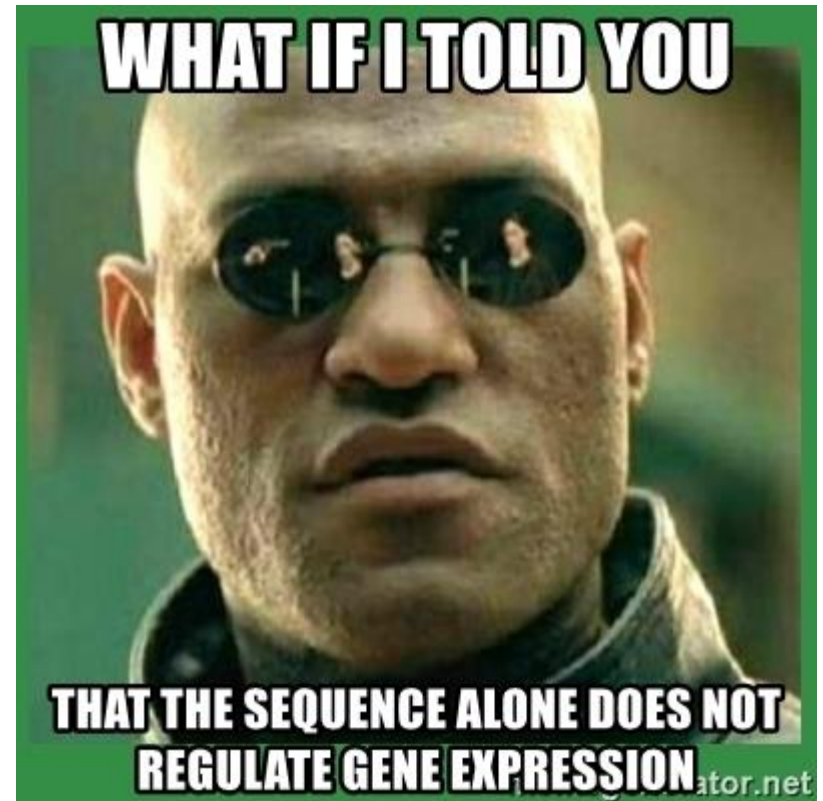
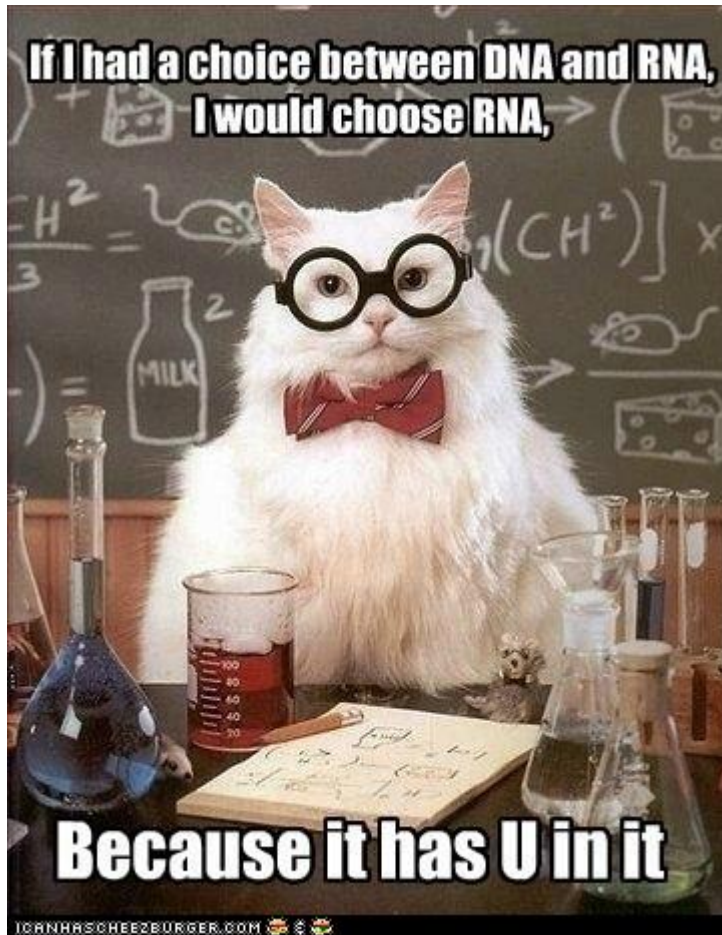
- RNA-sequencing and microarrays are generally used for high-throughput gene expression data, with the former eclipsing the latter.
- Pre-processing of RNA-seq data requires alignment of reads, transcript identification and quantification.
- RNA-seq allows us to capture information about mRNAs but also other RNA molecules.



# Additional Take-aways

- RNA-seq allows us to capture information about mRNAs but also other RNA molecules.
- For a beginner, it is easier to learn how the pipelines work by starting with bulk RNA-seq data.
- Innovative methods exploit the advantages of gene expression datasets – this is a space for amazing science that you could even do from home!
- Bioinformatics-driven science is a reality.
- Opportunities exist to use publicly accessible datasets to act on your curiosity!

**So, what are you curious about?**



Please check out Part 2 and see you next week for our online lecture!  
AIBN (Building 75) Level 4 West  
Email: [j.mar@uq.edu.au](mailto:j.mar@uq.edu.au)