# Lecture 25.2

# "Big data" in genomics experiments

Nicola Romanò
nicola.romano@ed.ac.uk

浙江大学爱丁堡大学联合学院
ZJU-UoE Institute

# Learning objectives

At the end of this you should be able to:

- Define "big data" in the context of biomedical research

- Discuss the issues associated with big data

- Describe a typical pipeline of genomic sequencing

# Big data – an historical perspective

- What is "big data"?

  *"Big data is a simple way of referring to data sets whose size grows beyond the ability of software and hardware tools to manage, capture, and process in a reasonable timeframe. Big data is where the amount of data is so massive that it becomes very difficult to control."*

  Iannarelli and O'Shaughnessy, 2015

- Term (in its current form) popularized in 2005 by Roger Mougalas from O'Reilly Media…

- … but the problem has been around much longer.

# Big data – an historical perspective



Early "big data" management in the 1930s – source Wired

- In 1937 the US government enacted the Social Security Act

- Data from 26 million Americans >3 million employers

- IBM developed punch card readers specifically for the job

# Big data and computing technology

- The sheer volume of data produced today requires

  - A lot of space (virtual  and physical) to **store** data
  - Large computing power to **analyse** data
  - Enhanced security to **protect** data



One of Google's data centres – Source Google

# Big data in biology

Biology has been affected as well

From "traditional", hypothesis-driven experimental strategies to large scale exploratory studies

# Deduction vs induction

# Big data in biology

Biology has been affected as well

From "traditional", hypothesis-driven experimental strategies to large scale exploratory studies

"Big data" is found in many domains of biomedical research

- Genomics

- Proteomics

- Imaging

- ...

# New technologies in genetics

Enormous efforts to sequence the human genome contributed to a big technological push in technologies related to genetics.



PCR – source: Wikipedia



Venter, J. C. *et al*. **Science** 291, 304–1351 (2001).

International Human Genome Sequencing Consortium, **Nature** 409, 860–921 (2001).

1985 – Kary Mullis develops PCR

2001 – The human genome is sequenced

# New technologies in genetics

Introduction of next generation sequencing (NGS) made sequencing price drop faster then predicted. Today you can sequence a genome for <$1000.

# *Traditional* Sanger sequencing

Based on PCR using "terminator nucleotides" (analogues that can be incorporated in DNA, but cannot have another nucleotide attached)

# Next generation sequencing

Also called "massive parallel sequencing". Improved speed and efficiency, reduced costs.

In most cases based on using microfluidics platform

Generally reads small fragments of DNA or RNA (150-300 bp)
Each fragment is called a "read"

State-of-the art machines can sequence up to >30M reads/hour

Extract DNA → Fragment & tag → Amplify → Read → Analysis

# Illumina sequencing

In Illumina sequencing amplification is obtained by "bridge method"



DNA fragments    Primers

DNA strands are attached
to cell surface at one end

Ends are attached to surface
by complimentary primers

Enzymes create double strands

Denaturation forms two
separate DNA fragments

Repetition forms clusters
of identical strands

Source: atdbio

# Illumina sequencing

## Sequencing by synthesis



Source: atdbio

# Illumina sequencing

At each step a picture with the spots corresponding to each cluster is taken

# Bioinformatics analysis

The next step is to analyse the data

Data needs to be aligned to the genome and counted

QC → Alignment → Counts → QC & filtering → Further analyses

|  | S1 | S2 | S3 | S4 |
|---|---|---|---|---|
| G1 | 120 | 0 | 30 | 34 |
| G2 | 85 | 22 | 15 | 0 |
| G3 | 0 | 0 | 2 | 88 |
| G4 | 150 | 230 | 223 | 743 |
| … | … | … | … | … |
| G20000 | 45 | 23 | 55 | 11 |

# Bioinformatics analysis

**Visualization techniques**
Heatmaps and clustering, PCA - see Workshop 7!

# Bioinformatics analysis

Differential expression analysis

- Finds genes under-/over-expressed between two samples

- Considers fold change and adjusted p-value

- Useful to find targets for further analysis



Jiang, 2016

# Bioinformatics analysis

**Pathway analysis**

- Finds molecular pathways enriched in the sample

- Genes can be mapped to function e.g. through GO annotations



Schabort, 2016

# From bulk to single cell sequencing

Recent advances in NGS and microfluidics technology have allowed sequencing of the transcriptome of single cells.

**10X sequencing**



**Generally $10^3$-$10^5$ cells sequenced – $10^4$-$10^5$ reads/cell**

**Other methods (e.g. SmartSeq 2): ~$10^6$ reads/cell**

# scRNAseq analysis

| | Cell 1 | Cell 2 | Cell 3 | ... | Cell 12000 |
|---|---|---|---|---|---|
| G 1 | 19 | 0 | 150 | ... | 0 |
| G 2 | 0 | 0 | 11 | ... | 153 |
| G 3 | 0 | 0 | 22 | ... | 80 |
| ... | ... | ... | ... | ... | ... |
| G 13000 | 15 | 81 | 0 | ... | 22 |

Dimension reduction
(e.g. PCA, t-SNE, UMAP...)

| | Dim 1 | Dim 2 |
|---|---|---|
| C 1 | 2.34 | -0.35 |
| C 2 | 3.15 | 1.87 |
| ... | ... | ... |
| C 12000 | -2.43 | -0.88 |

Clustering



Fidanza et al., 2019

# t-SNE dimension reduction


Fidanza et al., 2019

**t-Distributed Stochastic Neighbour Embedding**

Developed in 2008 by developed by Laurens van der Maaten and Geoffrey Hinton

Particularly suited to high dimensional data, such as that of scRNAseq experiments.

Contrary to PCA is not a linear projection, so it can capture non-linear relationships.

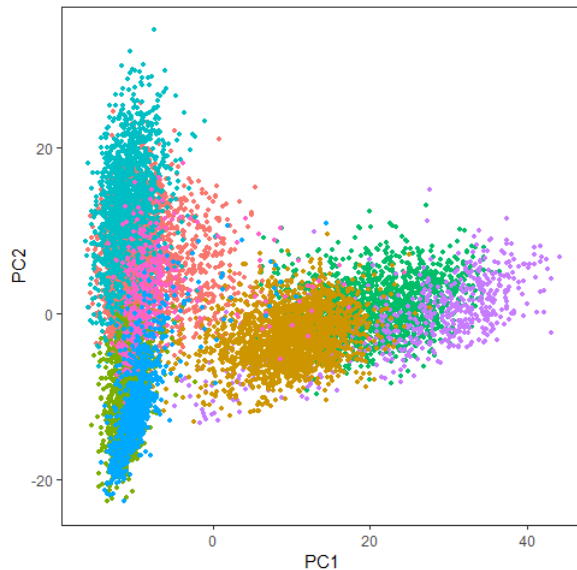Newer, non-linear methods are being constantly developed, such as UMAP (Uniform Manifold Approximation and Projection for dimension reduction).

# What are the issues?

An average file from a RNAseq experiment with ~50M reads is about 20 Gb

This has issues regarding

- Computing power

- Storage space

- Costs

- Data management plans are now required in most grants applications

- Bioinformatics skills are required (although pre-existing analysis pipelines make things a lot easier than they used to be)

# Not only sequencing...

Other sources of big data exist including

- Proteomics (and single-cell proteomics) experiments where MS is used to characterise the proteome

- 4D imaging experiments, high content screening, super-resolution methods

- Many more will be coming!

# Conclusion

New technologies are changing the way we do biology

Large amount of data are now easy and cheap to generate…

…at the cost of not being able to understand what those data mean!

New skills and technologies will be needed ever more in the future!

- Data analysis
- Machine learning
- Data visualization
- …