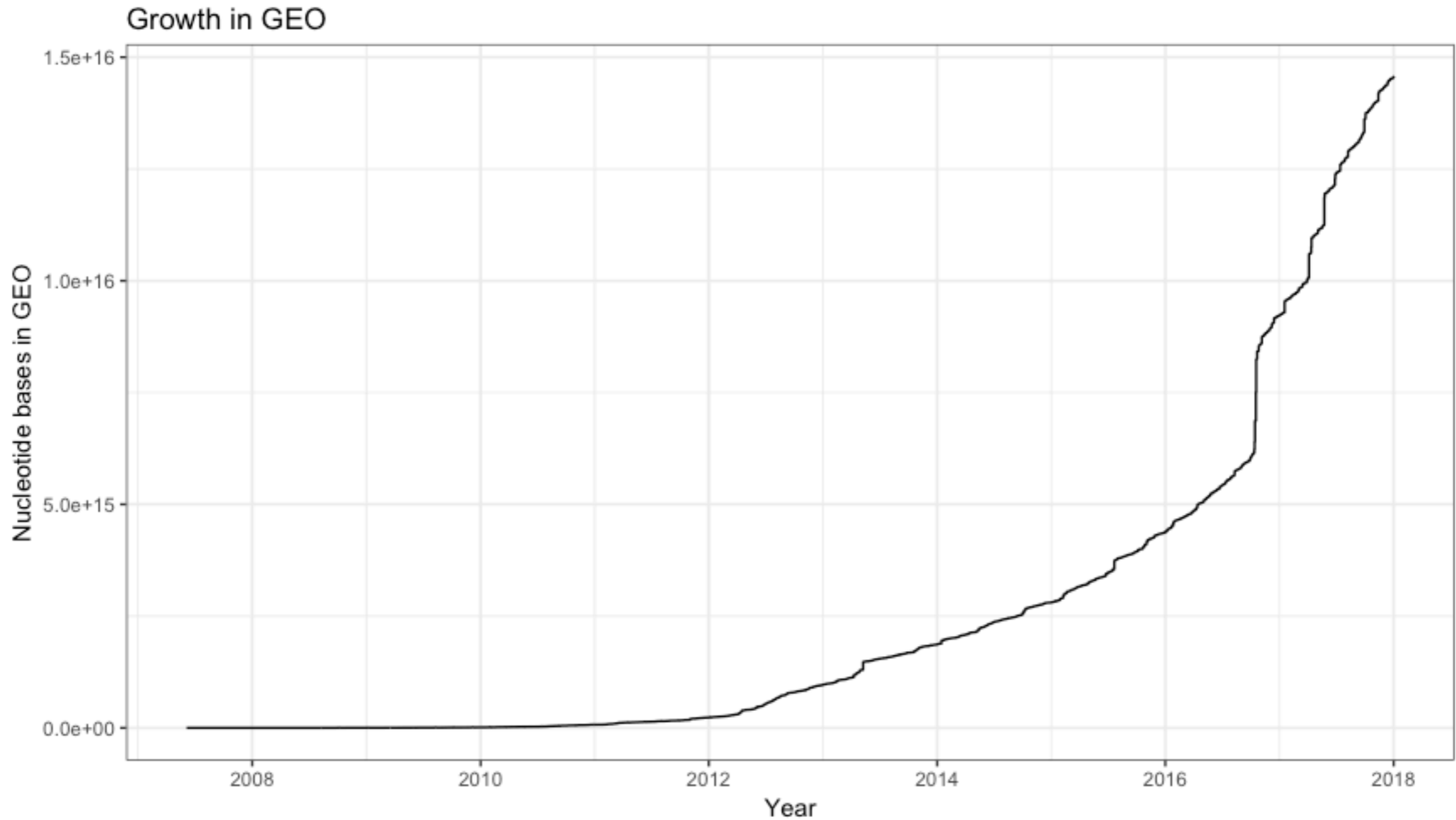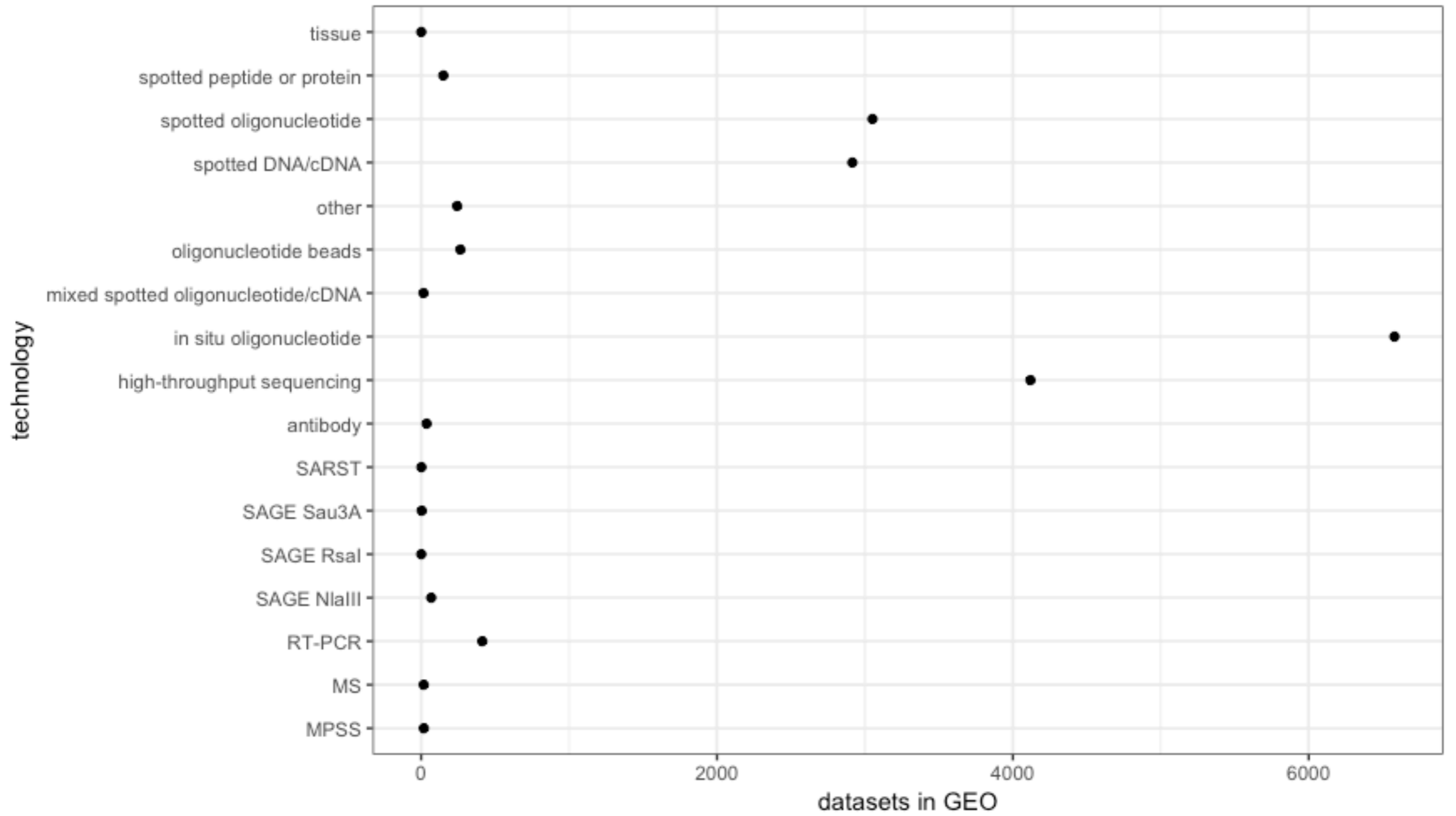# Data Integration in the Era of "Omics": Current and Future Challenges

Sandelin Lab
The Bioinformatics Centre, Dept. of Biology
Biotech Research and Innovation Center
(BRIC),
University of Copenhagen

UNIVERSITY OF COPENHAGEN

BRIC
Biotech Research &
Innovation Centre

# The classical "why you should get a degree in bioinformatics" slide



Growth in GEO

There is a reason that we give away our data and put it in a database - it is so that other people can use it. Combining data might produce insights that would not be possible with one data set alone.

## There are many good reasons to integrate data…

There are large sets of data available, for free.

Re-using data is cheap and efficient: most data is only surface analysed - maybe you can find new things?

Conversely: cheap to produce very rich, new data

The most interesting data analysis is when you can combine two or more data sets

Same experiment from many individuals
Same experiment from many tissues/cells
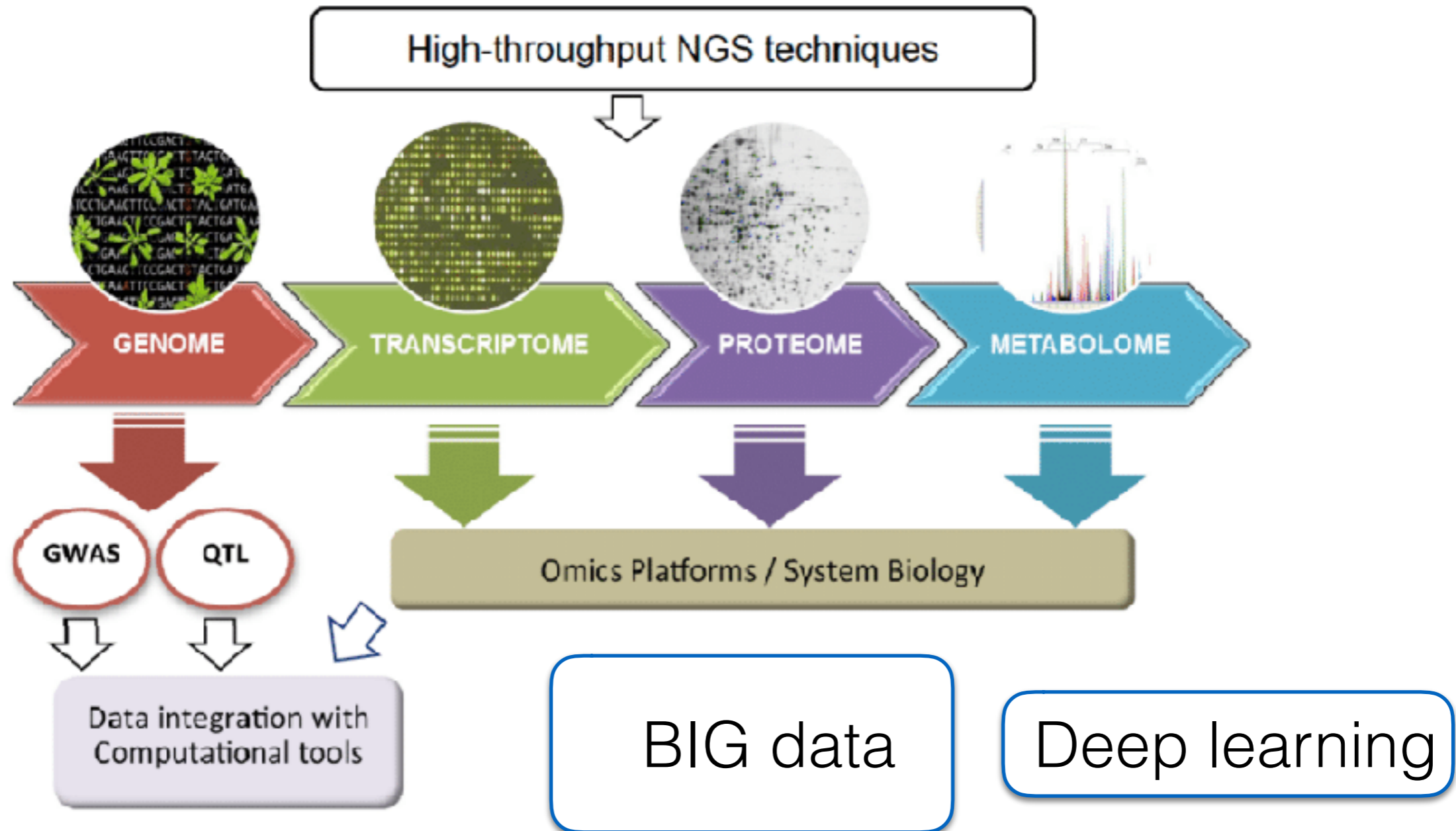Different experiments on same individual
Different experiments from many tissues/cells

3 scenarios:
- Design your own large study
- Take data from other big study (or possibly add to it)
- Take data from multiple studies (single or large studies)

But how to do this? The typical image around this is to make impressive arrows to a box called "Data integration" or "Systems Biology"
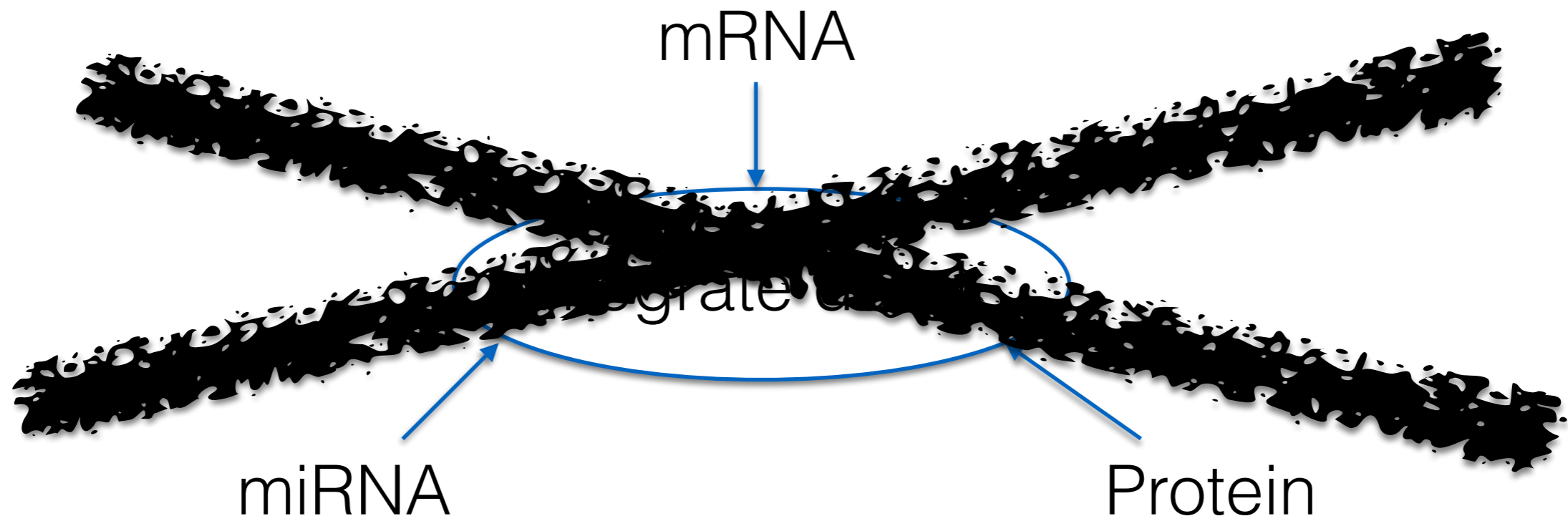


BIG data

Deep learning

## Lesson 1: Start with the question(s), not the data

What is the goal of your data integration? What type of questions do you want to answer?

Then: what data do you need, and how do you imagine these things should be analysed?

Example: I want to find out miRNA regulation of gene expression and protein expression, across individuals within same cell type

mRNA

miRNA

Protein

**Analysis (example):**
- Find genes with sites for each miRNA
- For each such gene, correlate mRNA expression, protein expression with of miRNA expression across samples.
- Hope to find convincing negative correlations

# Data needed

| **miRNA expression** | | **mRNA expression** | | **Protein expression** |
|---|---|---|---|---|
| Patient A | → | Patient A | → | Patient A |
| Patient B | → | Patient B | → | Patient B |
| Patient C | → | Patient C | → | Patient C |
| Patient D | → | Patient D | → | Patient D |
| Patient E… | → | Patient E… | → | Patient E… |

…from same samples, and same gene models

## Lesson 2.1 If you don't have comparable data, don't integrate!

| miRNA expression | mRNA expression | Protein expression |
|---|---|---|
| Patient A → | Patient A → | Patient A |
| Patient B → | Patient B → | Patient B |
| Patient C → | Patient C → | Patient C |
| Patient D → | Patient D → | Patient D |
| Patient E… → | Patient E… → | Patient E… |

| miRNA expression | mRNA expression | Protein expression |
|---|---|---|
| Patient A | Patient F | Patient H |
| Patient B | Patient G | Patient I |
| Patient C | Patient H | Patient J |
| Patient D | Patient I | Patient K |
| Patient E… | Patient J… | Patient L… |

## If analysis require linked samples, anything less than linked samples wont make it

## ENCODE pilot example

- The ENCODE pilot aimed to test out many different technologies on just 1% on the genome, defined to capture both known and unknown loci (selected and random)

Problems:
- Each laboratory used their own technique, often their own (different ) cell lines. And often published their results in advance of the integrative study

- Led to large problems for the analysis persons, and for selling the whole project as a joint and merged data set

- Experiment type A shows B in a certain region  in cell C
- Experiment type D shows E in same region in cell G

Indecipherable

- It meant the integration just became cumulative: so and so many bases have evidence for experiments of this and this type, and there was very limited narrative. Made it very hard to publish - was envisioned as 5 Nature papers but ended as 1.

## Lesson 2.2 . Beware of batch effects and confounders - leads to non-comparability

All data contains signal and noise.

Noise can be random or systematic. Systematic noise is hugely effected by experimental design: especially batch effects and confounding effects
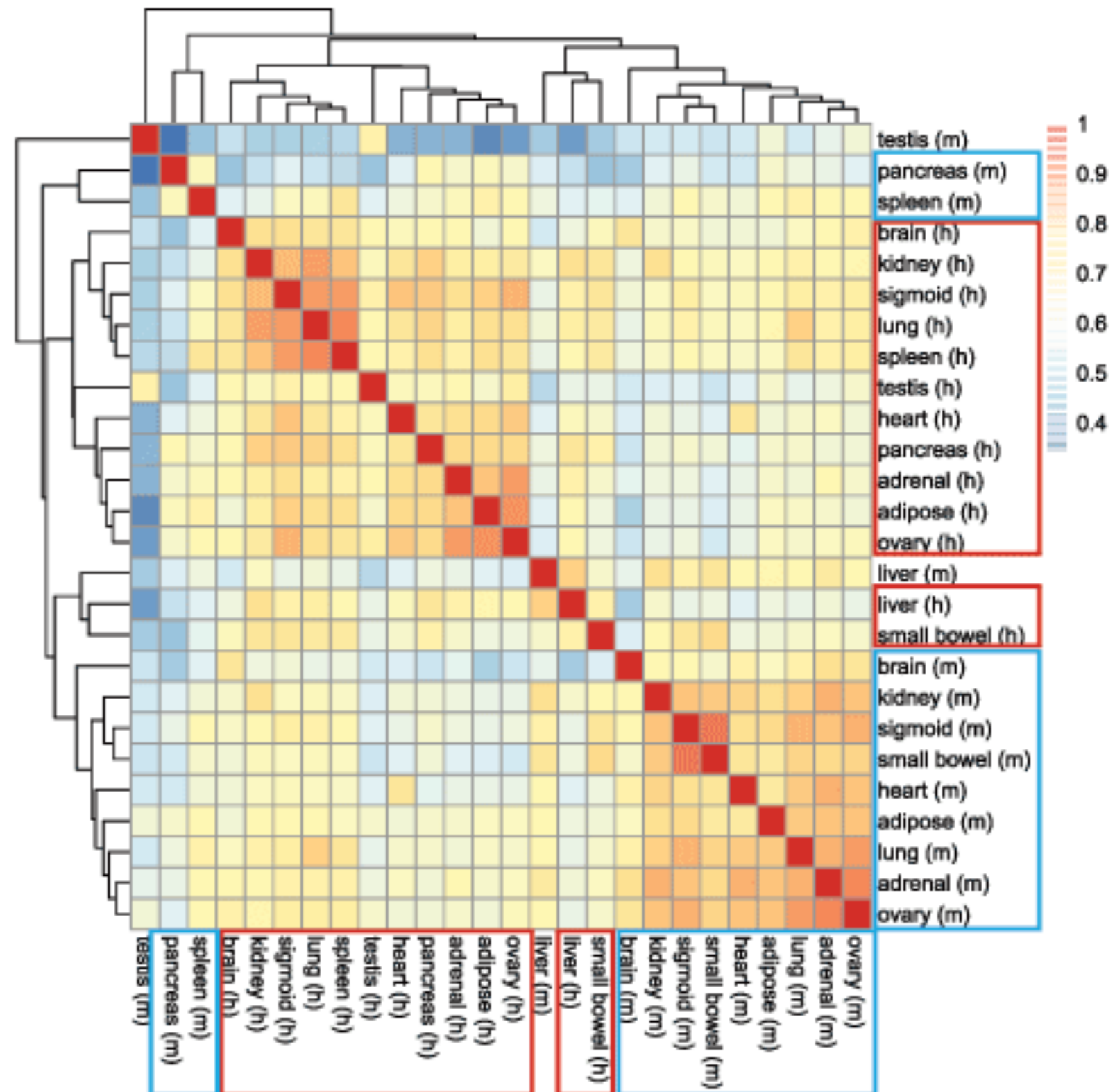
When dealing with large experiments, across multiple platforms, good experimental design becomes even more important - and sanity check afterwards.
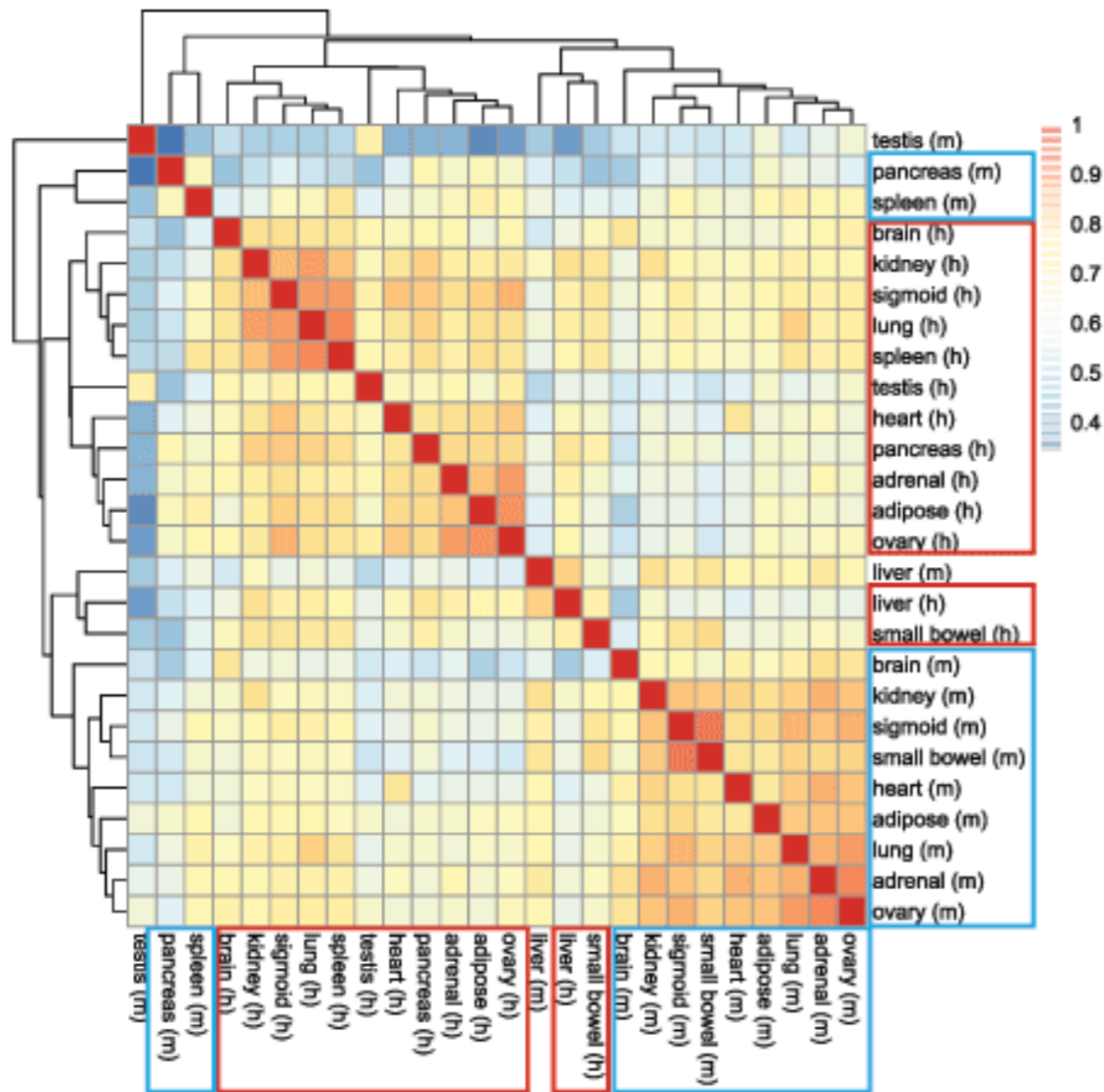
This is important as you may start reporting things that you are actually not interested in…

Analysis of gene expression across tissues , human vs mice:
Do tissue or species pair up?

Original study (Lin *et al.*: *PNAS.* 2014) found that all almost human tissue clustered together, so species is more important than tissue. This was very surprising as it went against many previous studies

Sequence study design (sequencer ID, run ID, lane number):

| D87PMJN1 (run 253, lane 7) | D87PMJN1 (run 253, lane 8) | D4LHBFN1 (run 276, lane 4) | MONK (run 312, lane 6) | HWI-ST373 (run 375, lane 7) |
|---|---|---|---|---|
| heart | adipose | adipose | heart | brain |
| kidney | adrenal | adrenal | kidney | pancreas |
| liver | sigmoid colon | sigmoid colon | liver | brain |
| small bowel | lung | lung | small bowel | spleen |
| spleen | ovary | ovary | testis | 🔴 human |
| testis | | pancreas | | 🔵 mouse |

Could this be due to bath effects? Gilad, 1000Research 2015, 4:121 tried to remove this effect using the COMBAT algorithm

# Non-batch-corrected

# Batch-corrected



Which one is correct?
Actually, we can never know.

## Lesson 3: Common notation, reference points and agreed data simplifications is critical

| **miRNA expression** | **mRNA expression** | **Protein expression** |
|---|---|---|
| Patient A ⟶ | Patient A ⟶ | Patient A |
| Patient B ⟶ | Patient B ⟶ | Patient B |
| Patient C ⟶ | Patient C ⟶ | Patient C |
| Patient D ⟶ | Patient D ⟶ | Patient D |
| Patient E… ⟶ | Patient E… ⟶ | Patient E… |

- How can we really compare miRNAs of Patient A to mRNA expression in Patient A? It requires an agreed standard on genes and how they are measured and defined

- Is splicing important? Should I consider different gene isoforms? Or do I simply want one value per "gene"?

- These definitions are in effect agreed simplifications or conceptualisations of the data, and will shape all analysis downstream

Three examples of large scale data consortial and their data integration efforts - not their analysis but rather their overall approach
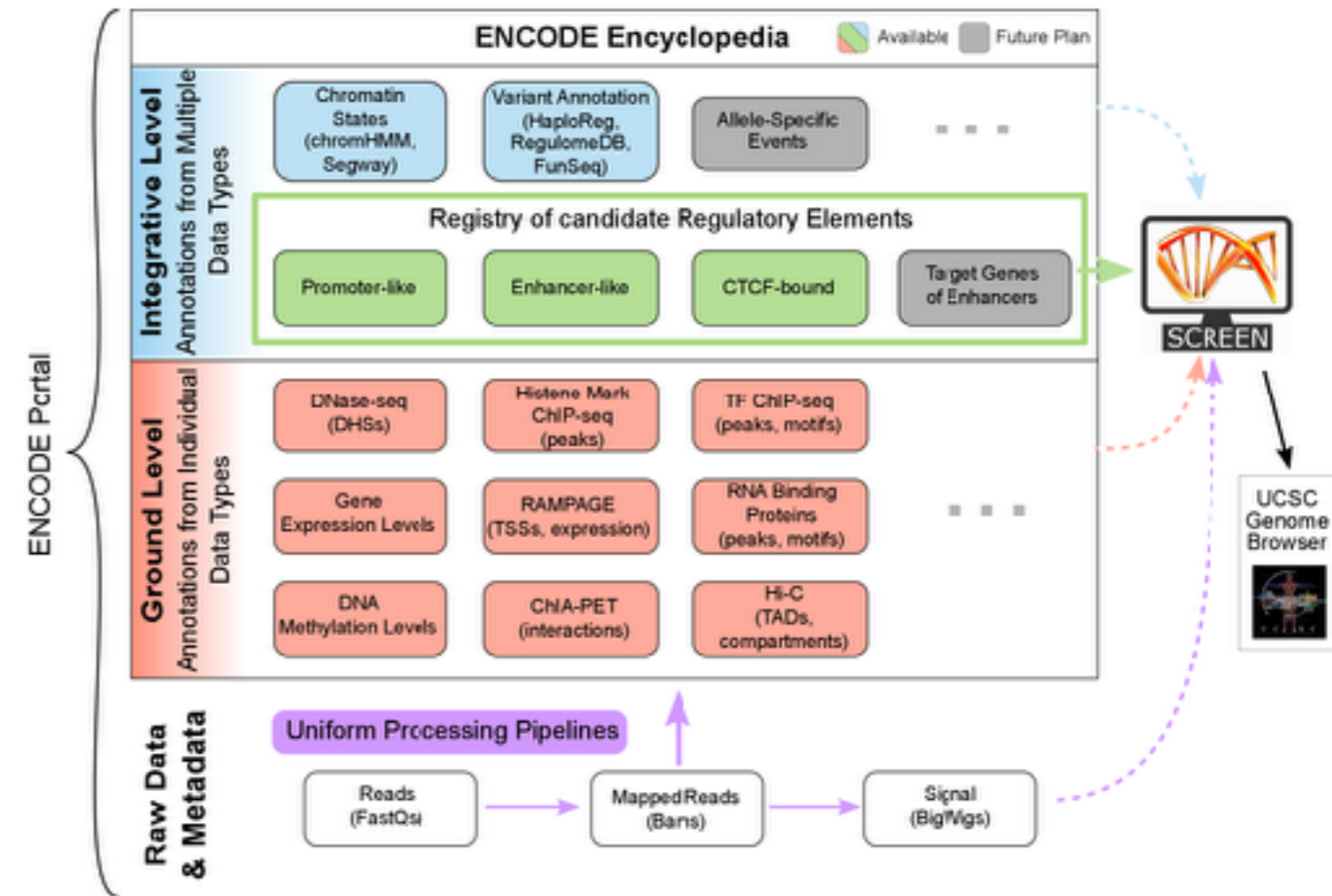
ENCODE
FANTOM
Cancer Genome Atlas

## Example 1: ENCODE phase 2

"Encyclopaedia" of genome elements (mostly cell lines) using a wide variety of techniques - mostly focused on gene regulation and gene expression, sometimes combined with knockdowns

Has a very clear strategy for data integration on multiple levels - starting from raw data, annotation from different data types and then integrated levels where relevant data sets are combined (combine many chromatin experiments into specific states, for instance)

**Common reference points**
1) For most experiments, the genome is the reference - peaks or blocks that then can be combined between experiments

2) GENCODE transcript models

Also has fantastic download section, and all data is "free before publishing"

From an data integration setting, ENCODE is worth studying since it is one of the largest project which started out very chaotic but then got to something that is perhaps the most mature

The ENCODE project also led to a lot of standards for eg ChIP data was established - extremely valuable

Possible downsides: very descriptive with somewhat unclear goals (narrative?), many labs contributing data (good or bad?)

## Example 2: FANTOM5 - atlas of transcription start sites

Question/Goal: Use CAGE-seq to profile transcription start site locations and usage across nearly all human cell types

Difference vs ENCODE: Only one technique, but a lot more samples. All data made at the same place, with many sample contributors, and many analysis spread over the world

Advantage: Much higher data control. On the other hand all data is locked down until a large paper is published

**Common reference points**

Promoters, or really clusters of CAGE tags across all experiments - a very large effort went into making these

A genome browser with all linked data

A wiki for collaborators with all raw and processed data
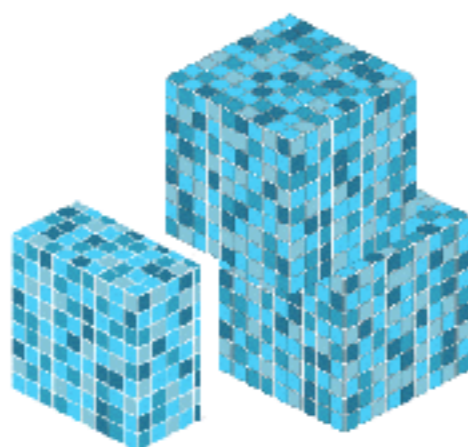
A data warehouse inc published
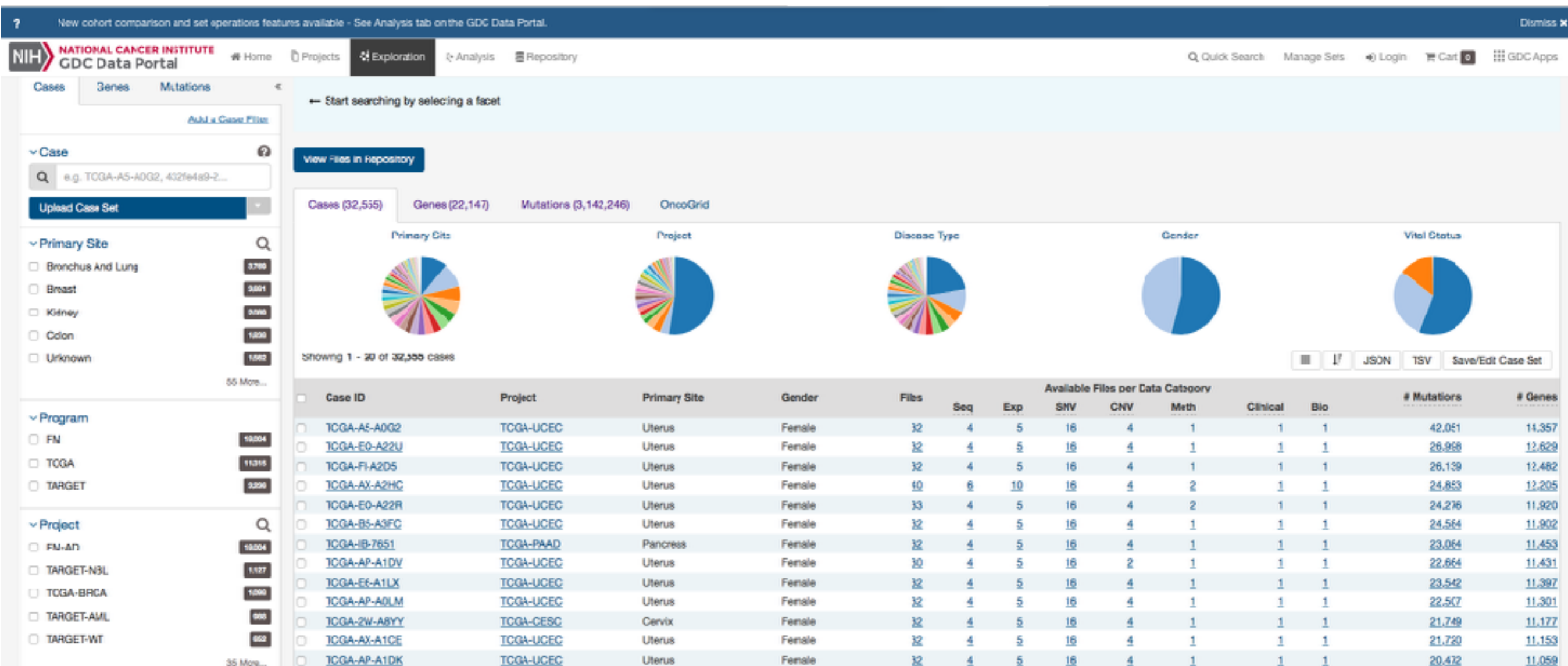
Example 3: Cancer genome atlas "TCGA"

# Much more complex reference points: on many levels: patients, cancer types tissue, genes, mutations



Fantastic overview on their strategy for different levels of data here:

https://cancergenome.nih.gov/abouttcga/aboutdata/datalevelstypes

## Summary

Data integration important, but often challenging - especially for ad-hoc data

No magic: Careful planning and setting goals, vision and narrative important. Once narrative is in place, integration "only" become a technical issue

Experimental design, data quality, outlier and bath identification is an important part of the integration - and it strangely enough becomes even more important if you collect other peoples data

Budget much more than you think to integration and integrative analysis - it takes time and good people

Consider what data that should be integrated: why should it be, what is the end goal. Avoid integration because of the "hype"

**There are data and projects where data integration makes zero sense.**

# Credits