

Extracting information from omics data

2pm Tuesday 9 April 2024; Data-driven approaches to understanding dementia

Charlotte Capitanchik, PhD

Postdoctoral Research Associate UK DRI @ KCL w/
Prof. Jernej Ule
Part-time Bioinformatics Lead @ Goodwright Ltd.

**(Gabriel) Mateus Bernado
Harrington, PhD**

Postdoctoral Research Associate UK DRI @ Cardiff
w/ Prof. Caleb Webber

KING'S
College
LONDON

CARDIFF
UNIVERSITY
PRIFYSGOL
CAERDYDD



**UK Dementia
Research Institute**

Outline

1. Overview of omics methods
2. Experimental design, Reproducibility
3. Pre-processing and pipelines
4. Data Visualisation
5. Multi-omic integration

All the Omics (well not all...)

1. Genomics

- a. Genotyping
- b. Whole genome seq

2. Transcriptomics

- a. Microarrays
- b. Bulk RNA-seq (short read or long read)
 - i. Random primed/3' end sequencing/5' end sequencing
- c. Single Cell/Nuclei
- d. Metabolic sequencing (e.g. SLAM-seq)
- e. RNA-protein binding (e.g. iCLIP, TRIBE)
- f. RNA structure (e.g. PARIS2)
- g. RNA modifications (e.g. miCLIP, from direct RNA seq, DART-seq)
- h. Ribosome profiling

3. Epigenomics

- a. DNA-Protein binding/histone modifications (ChIP-Seq, CUT&Tag, CUT&RUN)
- b. ATAC-Seq (Assay for Transposase-Accessible Chromatin using Sequencing)
- c. Bisulfite sequencing for m5C



All the Omics continued (still not all...)

1. Proteomics

- a. Tandem Mass Tag
- b. Targeted/Shotgun
- c. phosphoproteomics
- d. Protein-protein interactions (AP-MS)
- e. RNA-protein interactions (RBPome, OOPs)
- f. DNA-protein interactions...

2. Metabolomics

3. Lipidomics



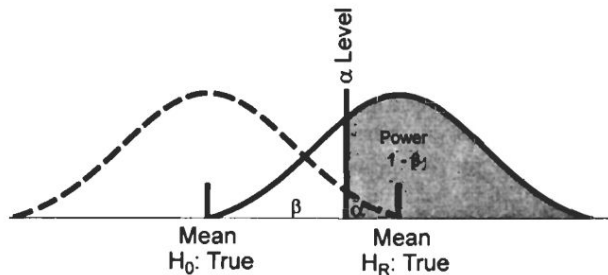
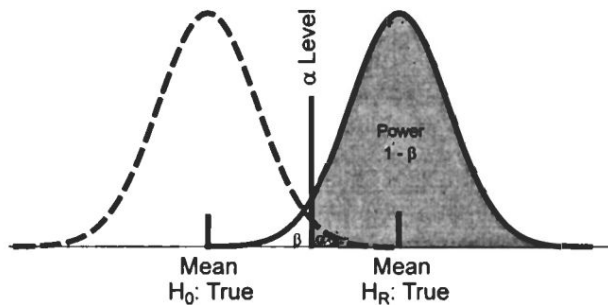
Experimental design - Controls

- This is a big topic worthy of careful consideration (and it's worth remembering you won't always have power to design the experiments you want to...), but one big point to highlight:
- **Controls!**
- Everyone always includes a negative control, but it's shockingly rare to see a positive control, use them if you can!

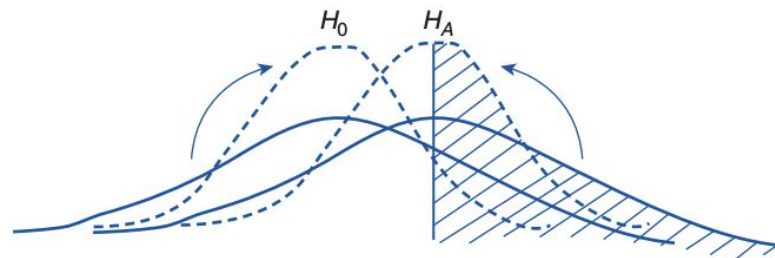


What impacts your ability to detect differences between conditions?

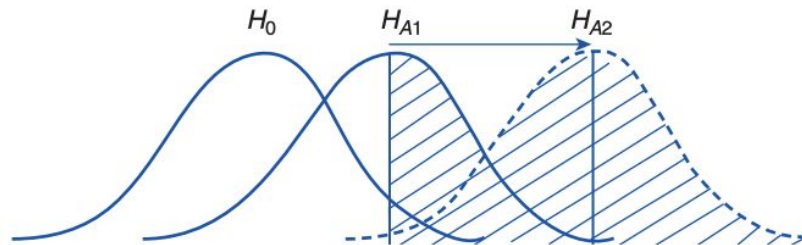
A.k.a. what variables affect statistical power?



↑ Variance
↓ Power



↑ Sample size
↑ Power



↑ Effect size
↑ Power

Practical implications

- If your effect size is **big**, you will need less replicates (samples) to observe it
- If you have a lot of variability between samples you can increase your power by making more replicates (beware 'N-hacking', but probs not that big a deal in omics? See Reinagel 2023, PLoS Biology.)
- In omics we can also think of read depth as “sample size” - eg. if there is a small difference between DEG (effect size) you will need higher read depth (sample size) to reliably detect it...
- People often talk of “signal-to-noise ratio” = if variance is high, power is reduced; the impact is worse if effect size is small.

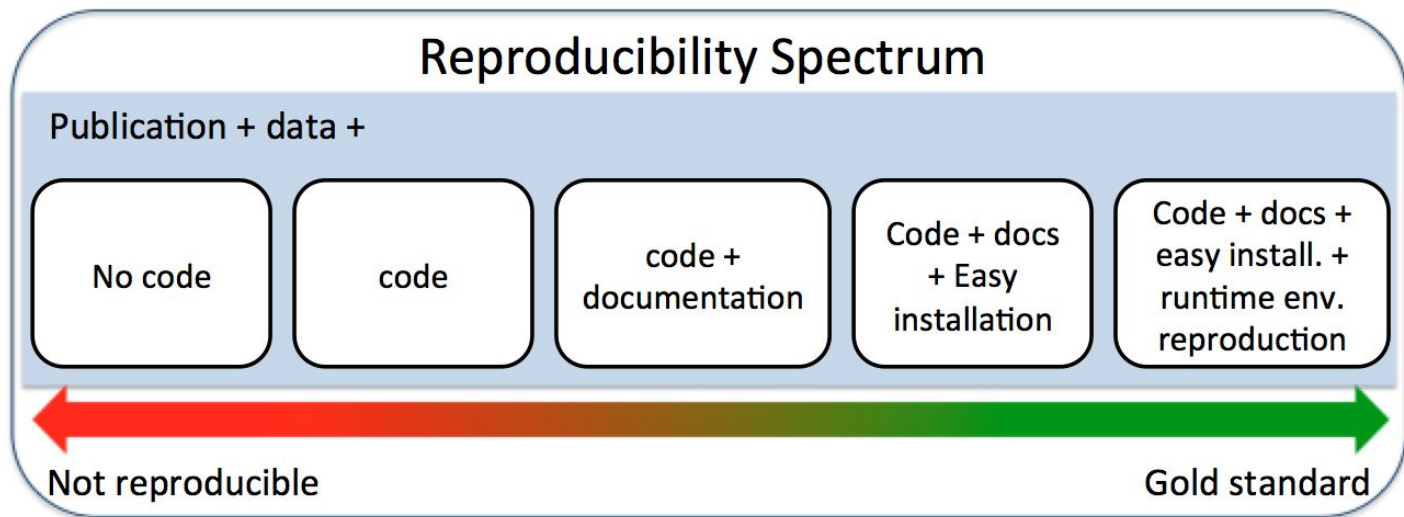
Reproducibility

- Reproducibility isn't just about lofty ideals of doing science that isn't a total load of rubbish
- If someone asked you to reproduce a figure you made 6 months ago from scratch, how easily could you do it?
- If you solved a niche problem 2 years ago and need to solve it again, how easily could you find and redeploy your solution?

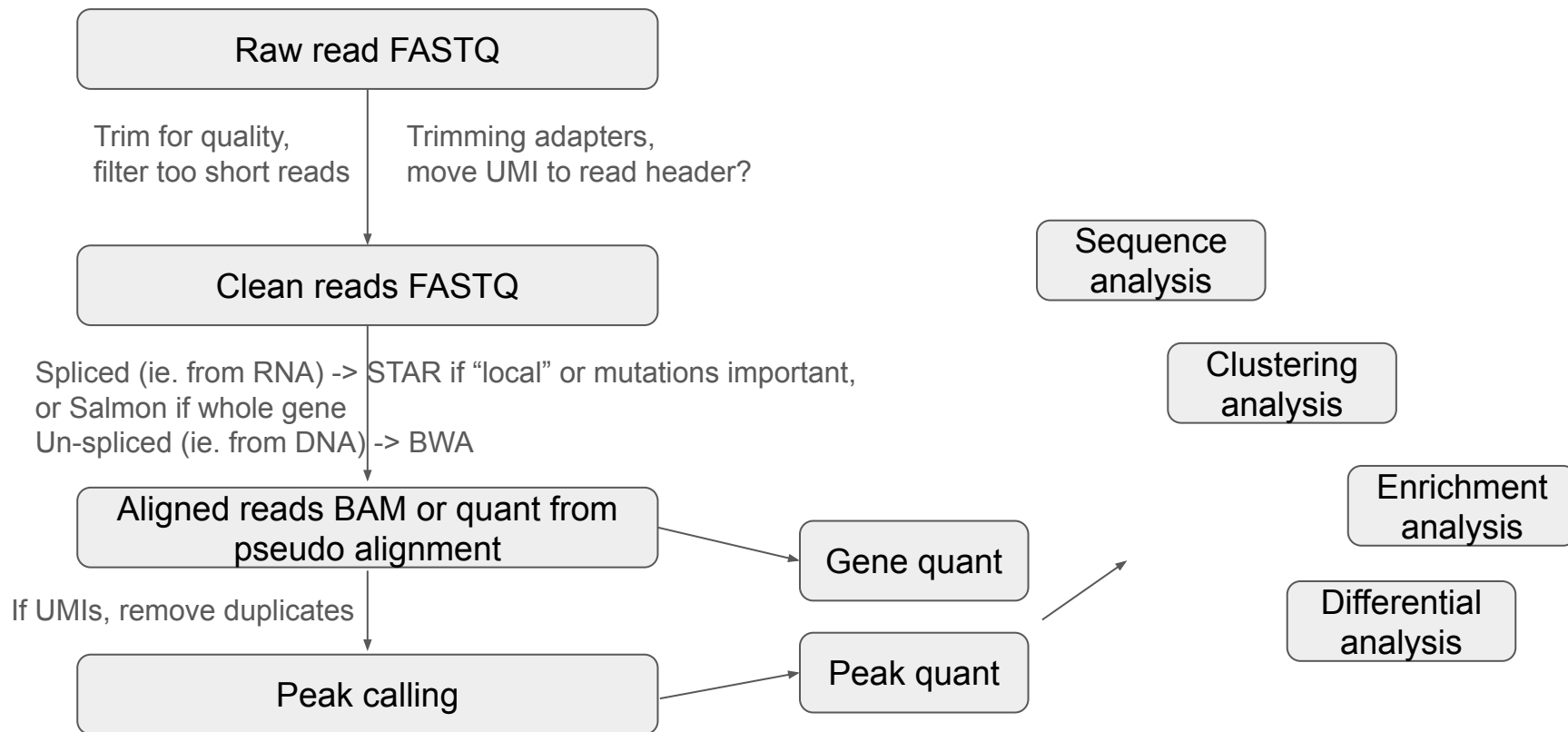


Reproducibility

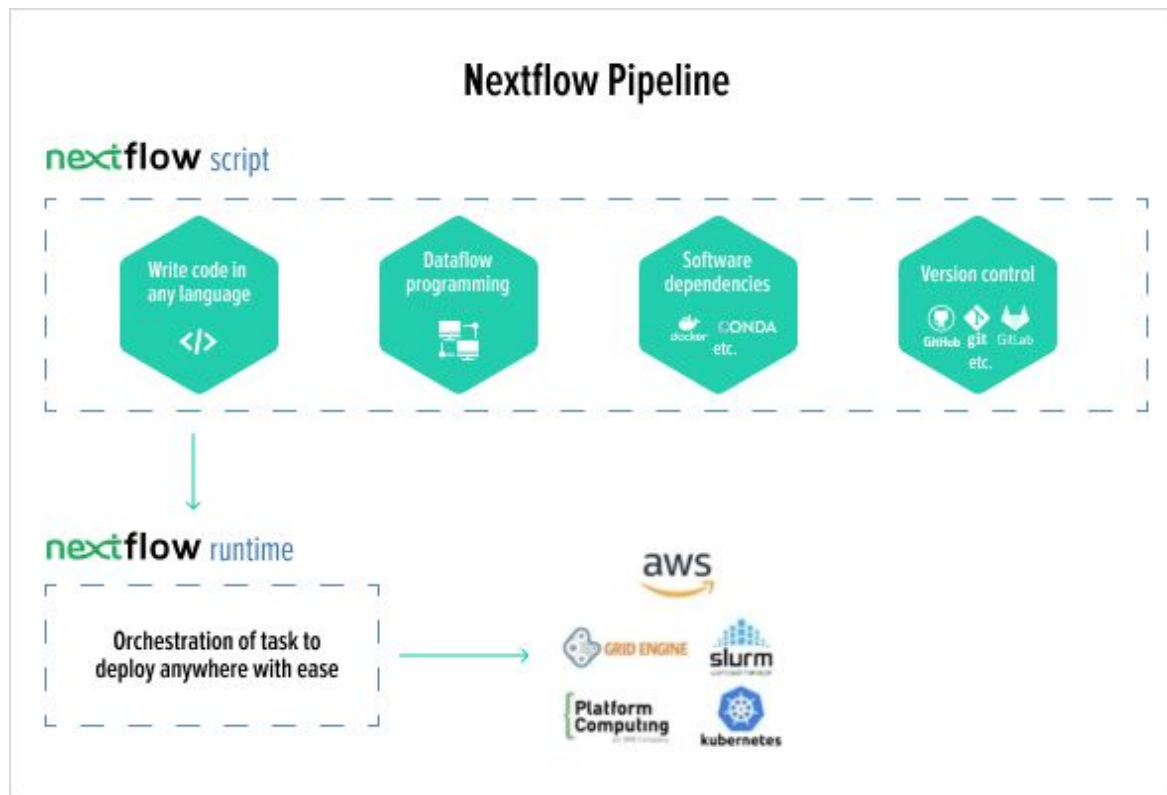
- Do future you a solid, and make sure your work can be reproduced, at least by you, ideally by others too!
- [The Turing Way](#) is a great guide to reproducibility if you want to learn more



Analysing xyz-Seq



Nextflow is a workflow language



Nf-core has lots of pipelines and modules ready to go



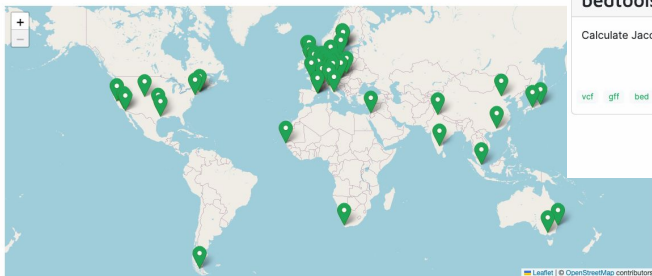
- 92 pipelines
- 1058 modules

Organisations

Some of the organisations running nf-core pipelines are listed below, along with a key person who you can contact for advice.

Note

Is your group missing? Please submit a pull request to add yourself! It's just a few lines in a simple YAML file.



nf-core/modules

Browse the 1201 modules that are currently available as part of nf-core.

bedtools

47 Name

bedtools_bamtobed

Converts a bam file to a bed12 file.

bam bed bedtools bamtobed converter

Included in: cutandrun radseq ssds

bedtools_closest

For each feature in A, finds the closest feature (upstream or downstream) in B.

bedtools closest bed vcf gff

bedtools_complement

Returns all intervals in a genome that are not covered by at least one interval in the input BED/GFF/VCF file.

bed gff vcf complement bedtools intervals

Included in: cutandrun

bedtools_coverage

computes both the depth and breadth of coverage of features in file B on the features in file A

bedtools coverage bam bed gff vcf histogram

Included in: radseq

bedtools_genomecov

Computes histograms (default), per-base reports (-d) and BEDGRAPH (-bg) summaries of feature coverage (e.g., aligned sequences) for a given genome.

bed bam genomecov bedtools histogram

Included in: cutandrun hicar nascent +1 more pipelines

bedtools_getfasta

extract sequences in a FASTA file based on intervals defined in a feature file.

bed fasta getfasta

Included in: readsimulator virarecon

bedtools_groupby

Groups features in a BED file by given column(s) and computes summary statistics for each group to another column.

bed groupby bedtools

bedtools_intersect

Allows one to screen for overlaps between two sets of genomic features.

bed intersect overlap

Included in: cirrna cutandrun nascent +1 more pipelines

bedtools_jaccard

Calculate Jaccard statistic b/w two feature files.

vcf gff bed jaccard intersection union statistics

bedtools_makewindows

Makes adjacent or sliding windows across a genome or BED file.

bed windows fai chunking

Included in: phasimpute radseq ssds

bedtools_map

Allows one to screen for overlaps between two sets of genomic features.

bed vcf gff map bedtools

bedtools_maskfasta

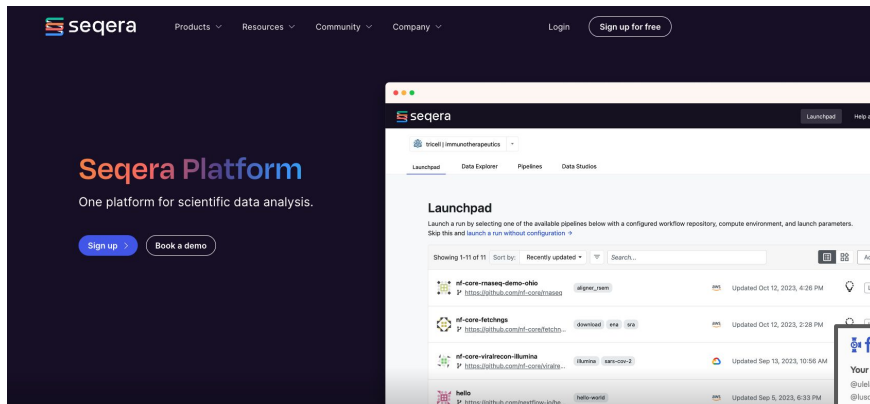
masks sequences in a FASTA file based on intervals defined in a feature file.

bed fasta maskfasta

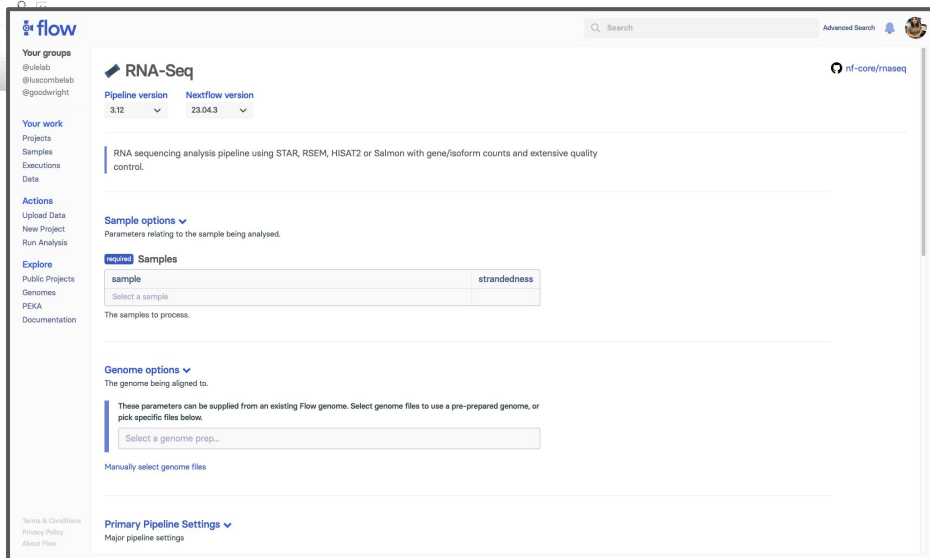
Included in: callingcards virarecon

Previous 1 2 Next

Nextflow pipelines can also be run through a GUI



app.flow.bio



A Nextflow-based bioinformatics analysis platform and open database



1 Upload your experimental data

2 Automatically store standardised data

3 Select parameters and run a pipeline in 1-click

4 Get your research insight and visualisations

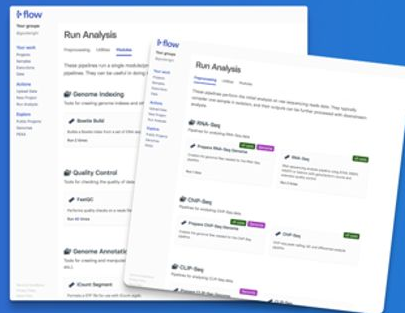
5 Add your results to a growing public database

Collate and view analysis history



View input/output files of Nextflow analyses in real time, along with curated data outputs

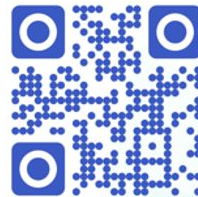
Run Nextflow pipelines down to individual modules (+ add custom pipelines)



View terminal history, stdout, and Nextflow logs



View analysis metadata and parameter history



Try now at flow.bio

No programming required - run Nextflow pipelines in the browser at the click of a button



goodwright
software for scientists



UK Dementia
Research Institute



Batch correction

- Omics data from brains is notoriously variable
- Data can be confounded by:
 - known variables, eg. time of day, time to autopsy, RIN, fixation, differences in dissection, laboratory, experimental batch, batch of library preparation
 - unknown variables eg. any of the above that weren't measured/recorded/communicated to you, GREMLINS
 - Dropout - eg. loss of rare cell types due to sample prep

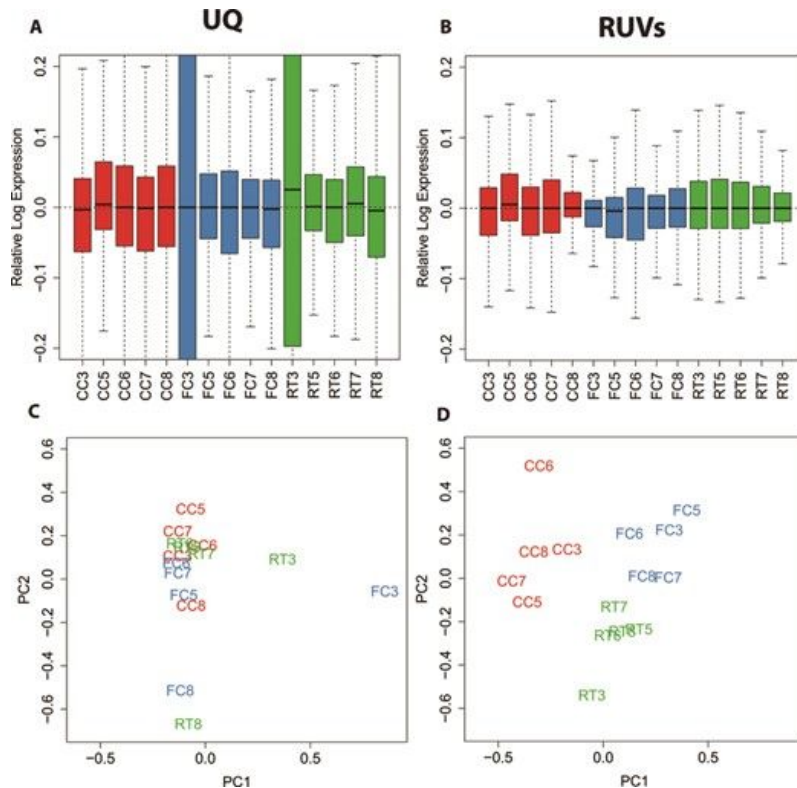
Batch correction

It is always best to try
and reduce variation
experimentally

BUT

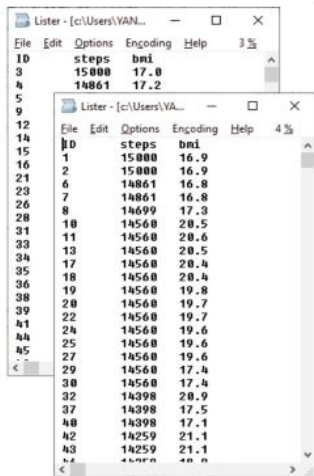
If you can't it is
possible to remove
known and unknown
confounders with

statistical modelling 🌈



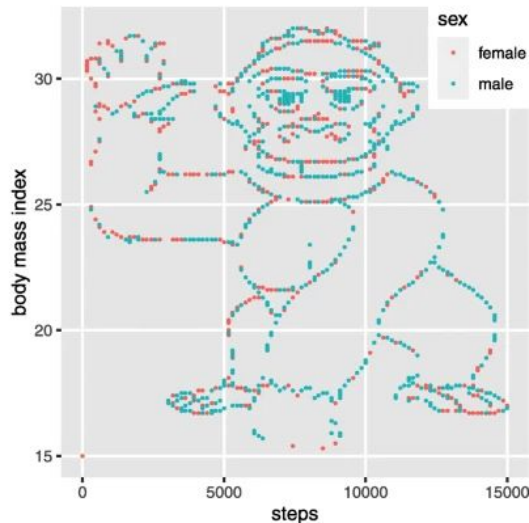
Always visualise your data

a



	steps	bmi
10	15000	17.0
3	14861	17.2
1	15000	16.9
2	15000	16.9
6	14861	16.8
7	14861	16.8
8	14699	17.3
10	14560	20.5
11	14560	20.6
13	14560	20.5
17	14560	20.4
18	14560	20.4
19	14560	19.8
20	14560	19.7
22	14560	19.7
24	14560	19.6
25	14560	19.6
27	14560	19.6
29	14560	17.4
30	14560	17.4
32	14398	20.9
37	14398	17.5
40	14398	17.1
42	14259	21.1
43	14259	21.1

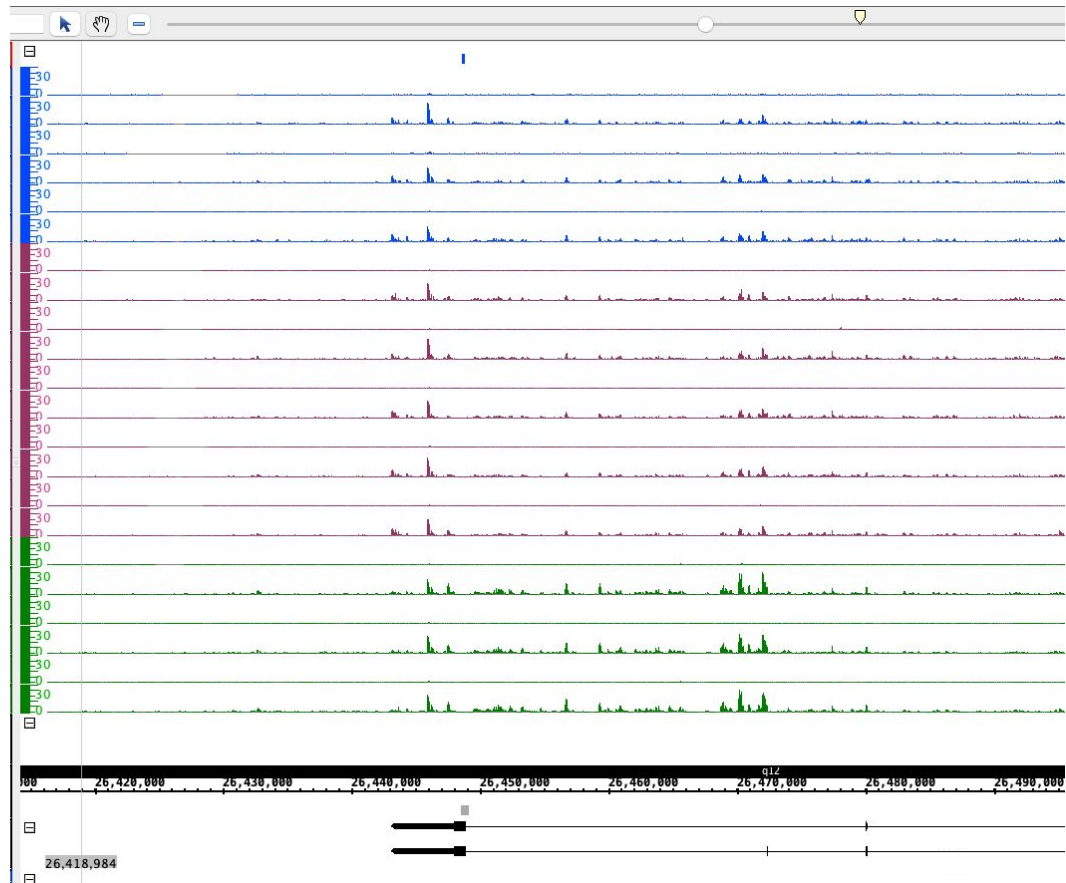
b



c

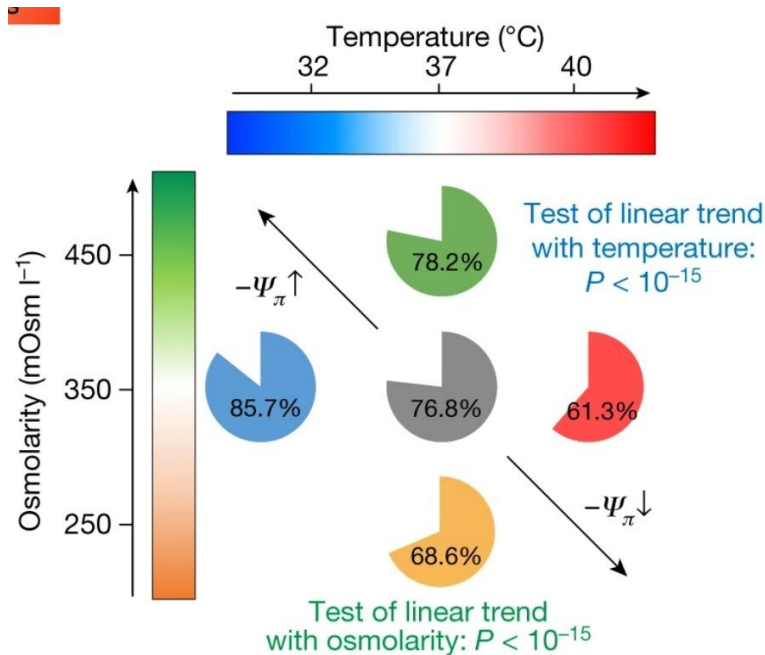
	Gorilla <u>not</u> discovered	Gorilla discovered
Hypothesis-focused	14	5
Hypothesis-free	5	9

Always visualise your data



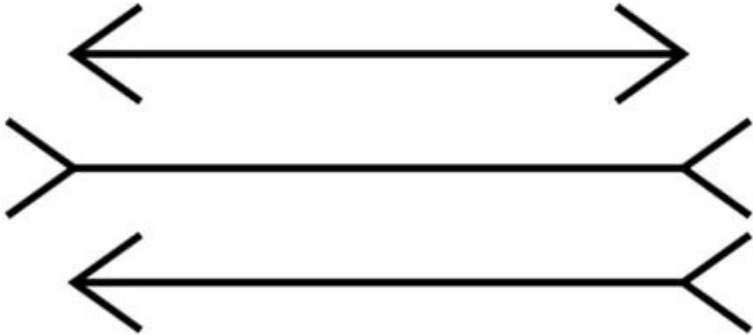
Responsible data visualisation

- With visualisation we aim to be: **clear, accurate & attractive** - in practise difficult to be all three!



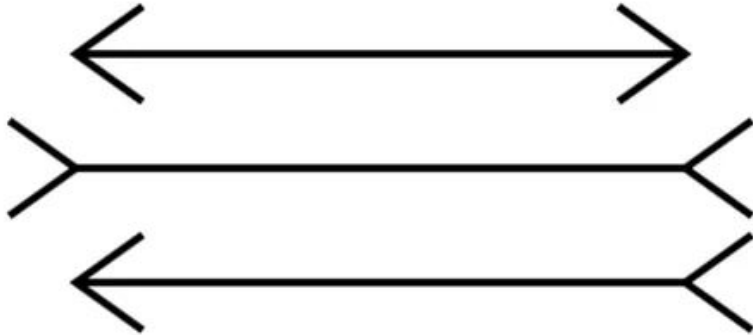
Responsible data visualisation

- With visualisation we aim to be: **clear, accurate & attractive** - in practise difficult to be all three!
- It's important to understand that visualisation is **rarely neutral**
- Humans have **biases in perception**, hence optical illusions!

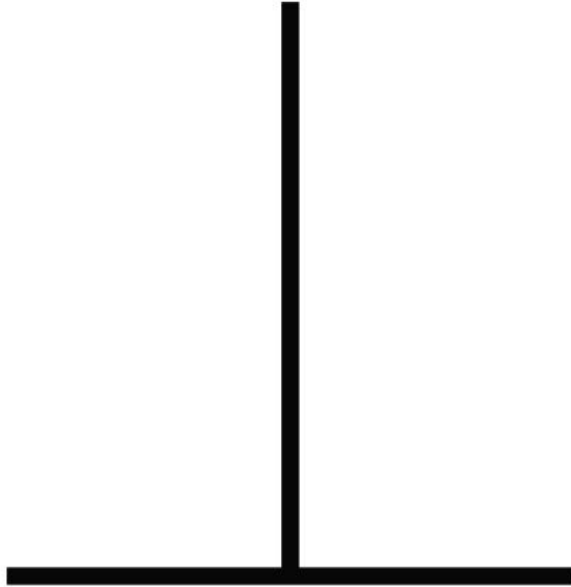


Responsible data visualisation

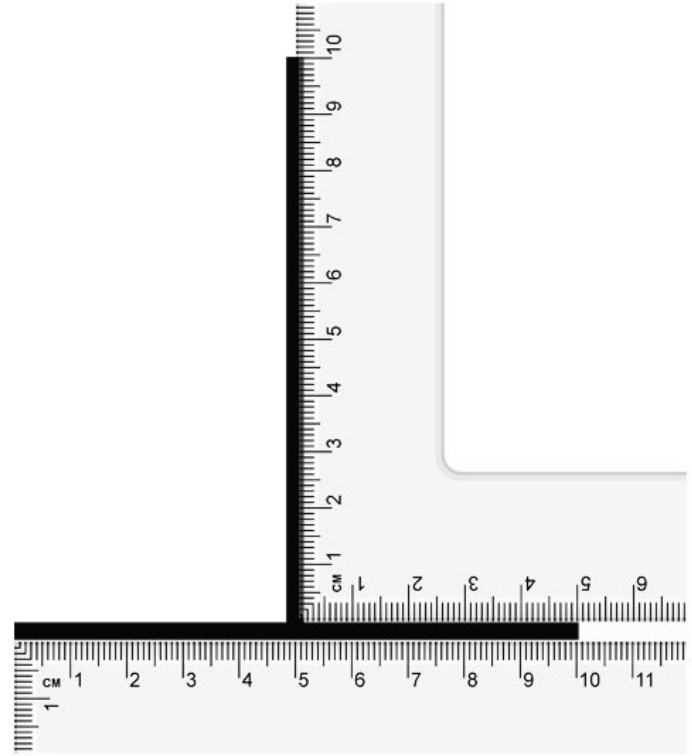
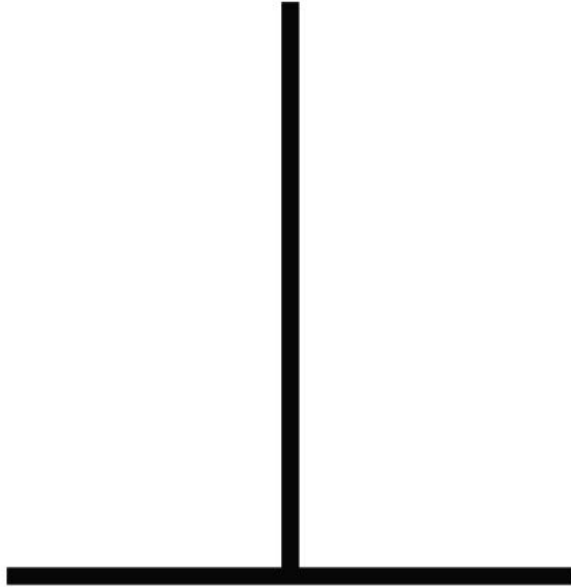
- With visualisation we aim to be: **clear, accurate & attractive** - in practise difficult to be all three!
- It's important to understand that visualisation is **rarely neutral**
- Humans have **biases in perception**, hence optical illusions!



Responsible data visualisation



Responsible data visualisation



Responsible data visualisation

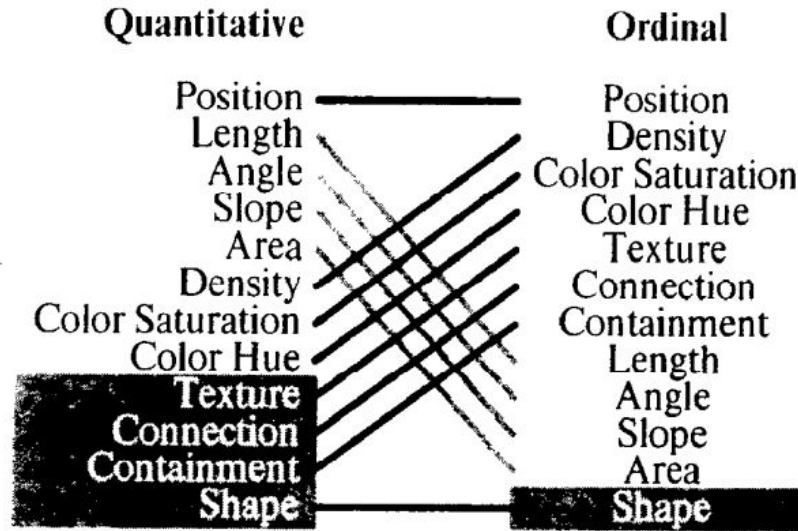
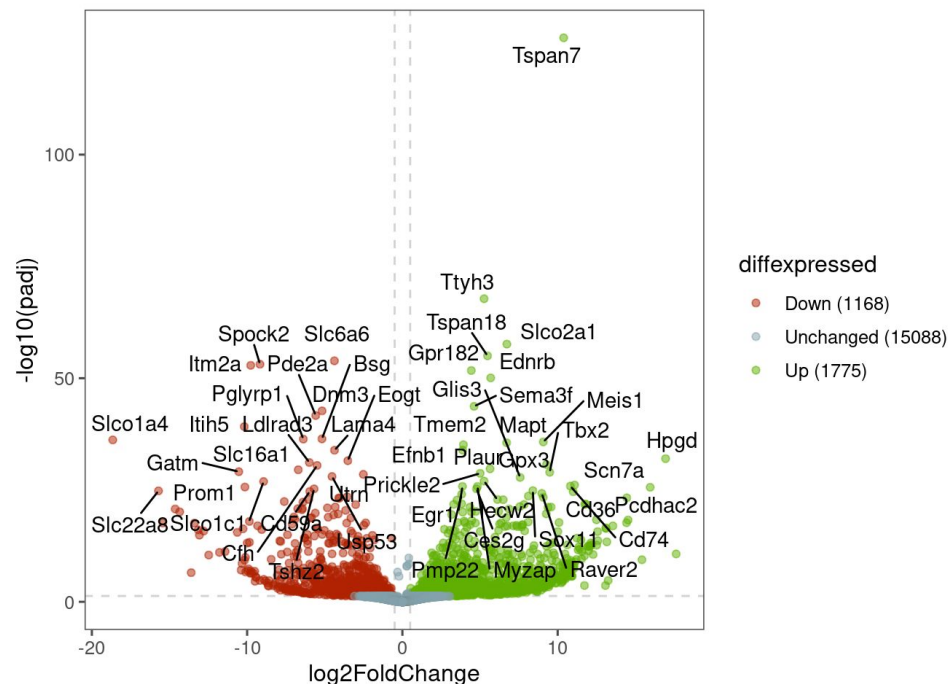
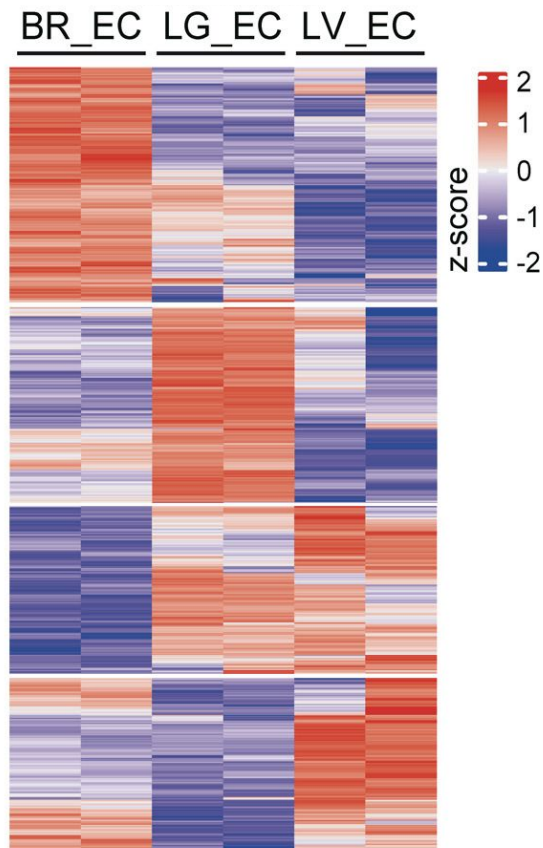


Figure 9. Ranking of perceptual tasks. The columns are for three different types of information. Tasks higher in the chart are perceived more accurately than tasks lower in the chart. The tasks shown in gray are not relevant to that type of information.

Heatmap vs. volcano plot for gene expression



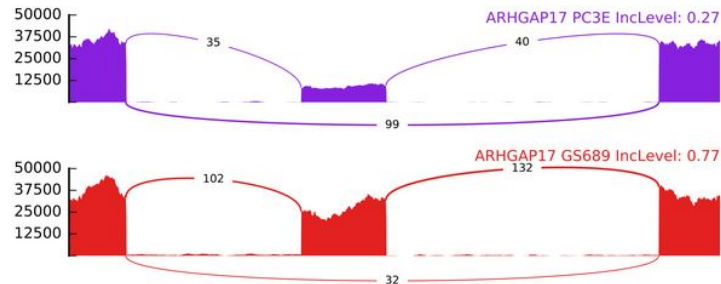
Multi-omic Data Integration

Ah yes, that thing I put in my grant...

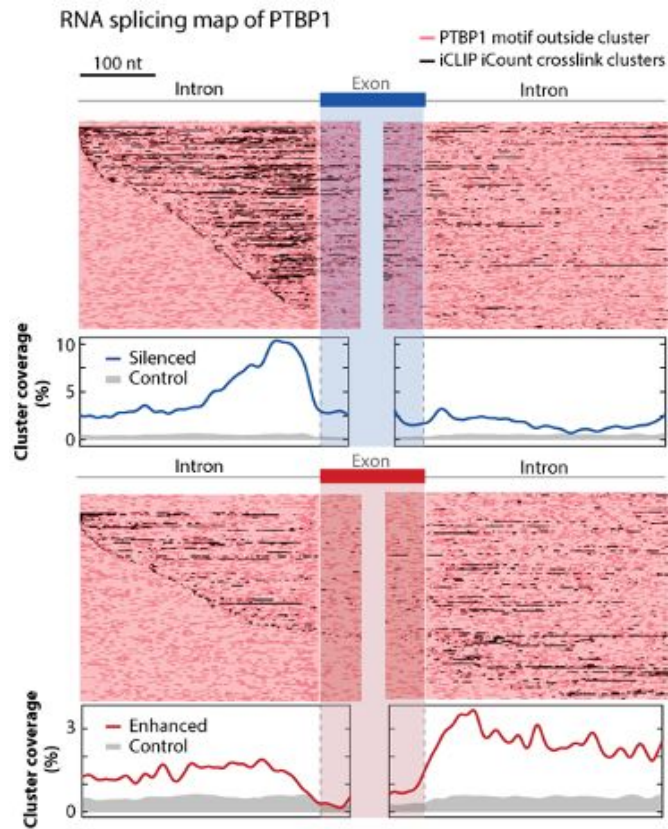
One example of integration of iCLIP and RNA-Seq...



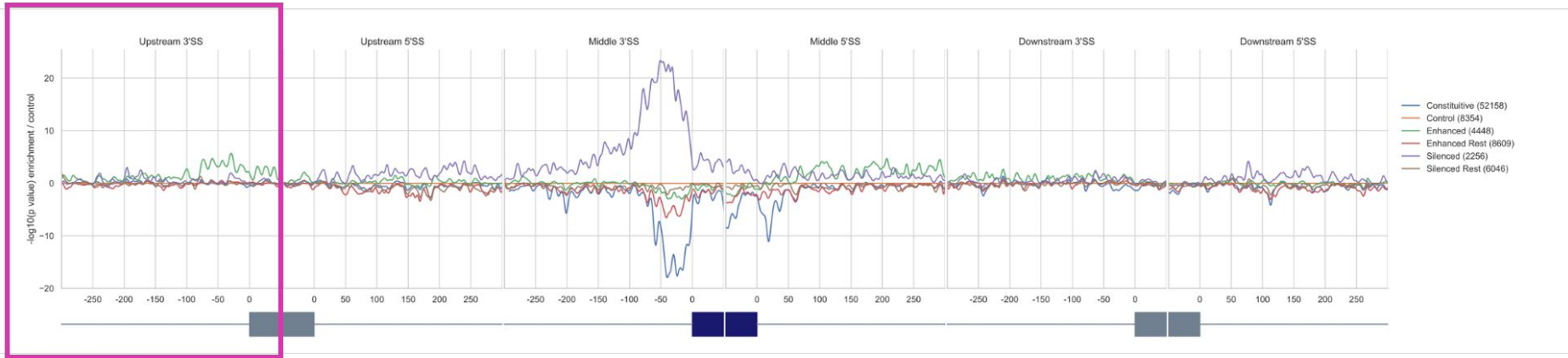
+



RNA splicing map



RNA splicing map

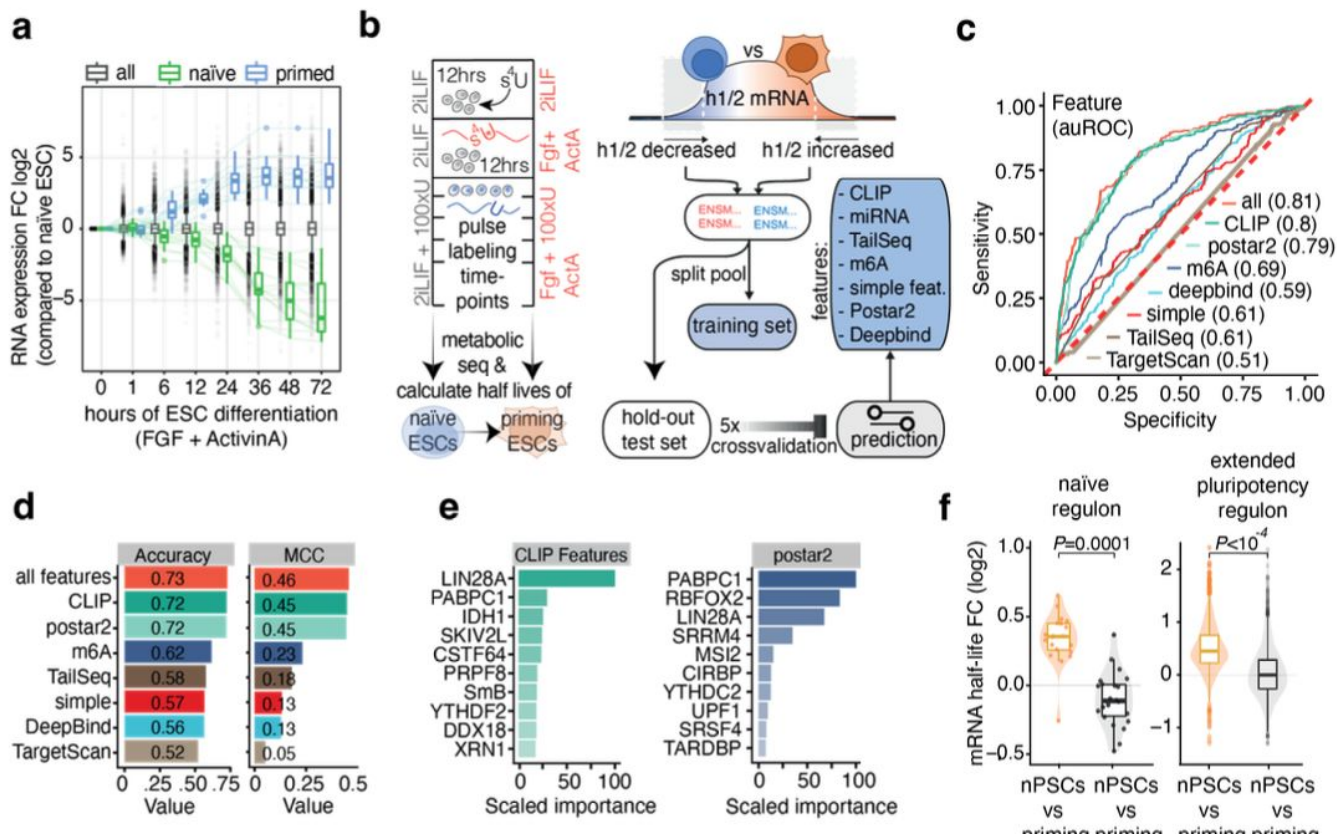


Ptbp1 RNA map

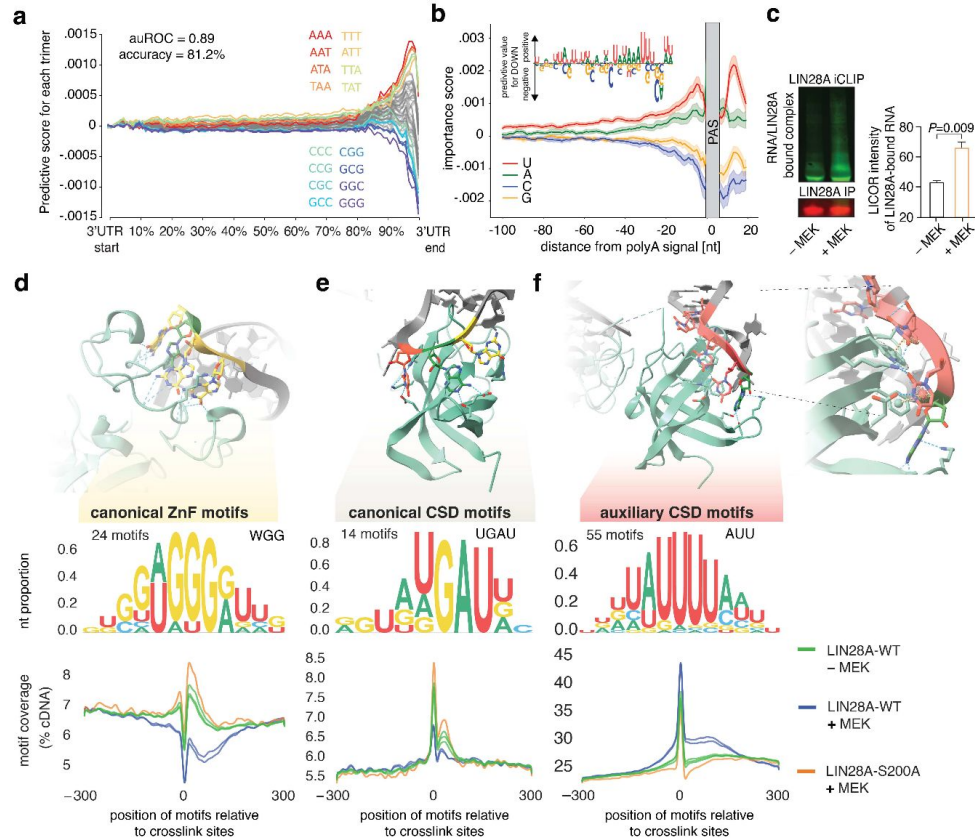
Modelling omics data

1. Chuck in all your data as features and try to predict something e.g. SVM, random forest, neural net..., can identify most predictive features.
2. Using only sequence as input can you predict something measured by omics data eg. gene expression, mRNA half-life, ribosome occupancy..., use backpropagation methods to figure out which sequences were important.
3. “Foundation models” training deep learning models to predict omics data tracks, use backpropagation methods to figure out which sequences were important or attention layer of transformer models.

Modelling omics data (1)



Modelling omics data (2)



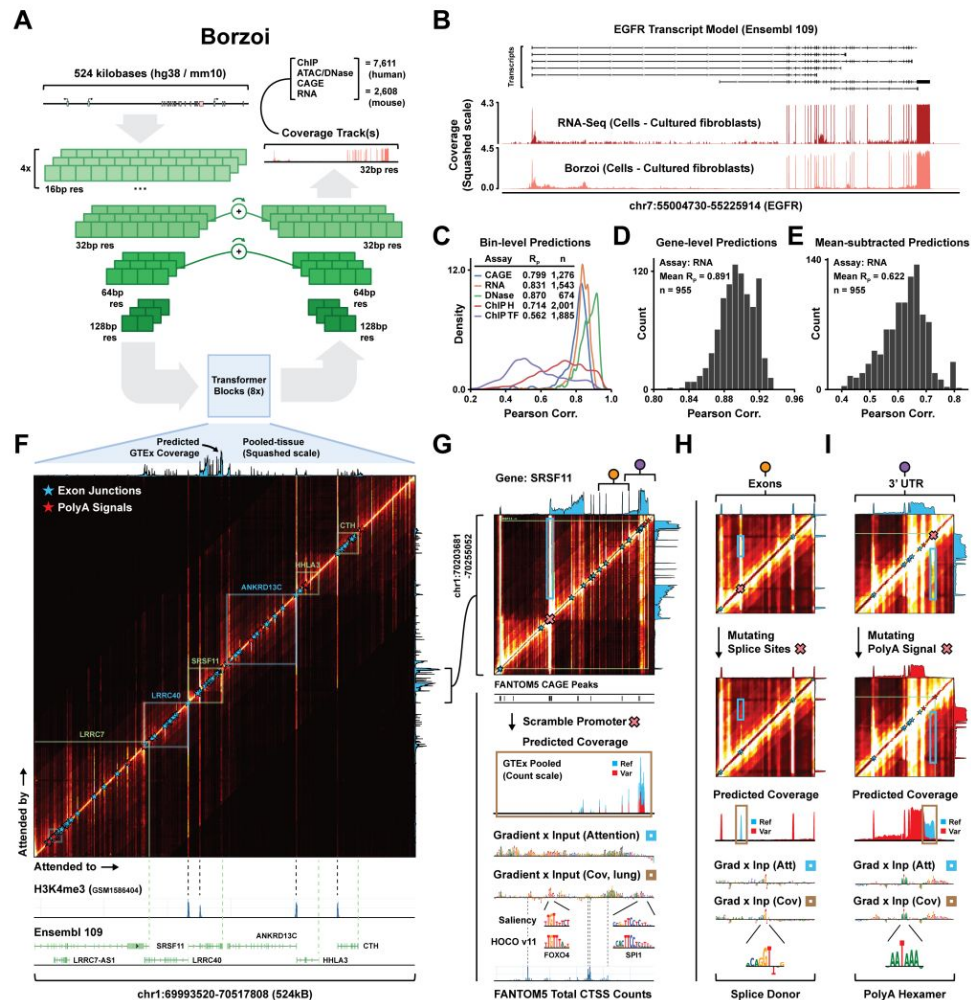
Modelling omics data (3)

Borzoï models are convolutional neural networks trained to predict RNA-seq coverage at 32bp resolution given 524kb input sequences

GTEx + ENCODE

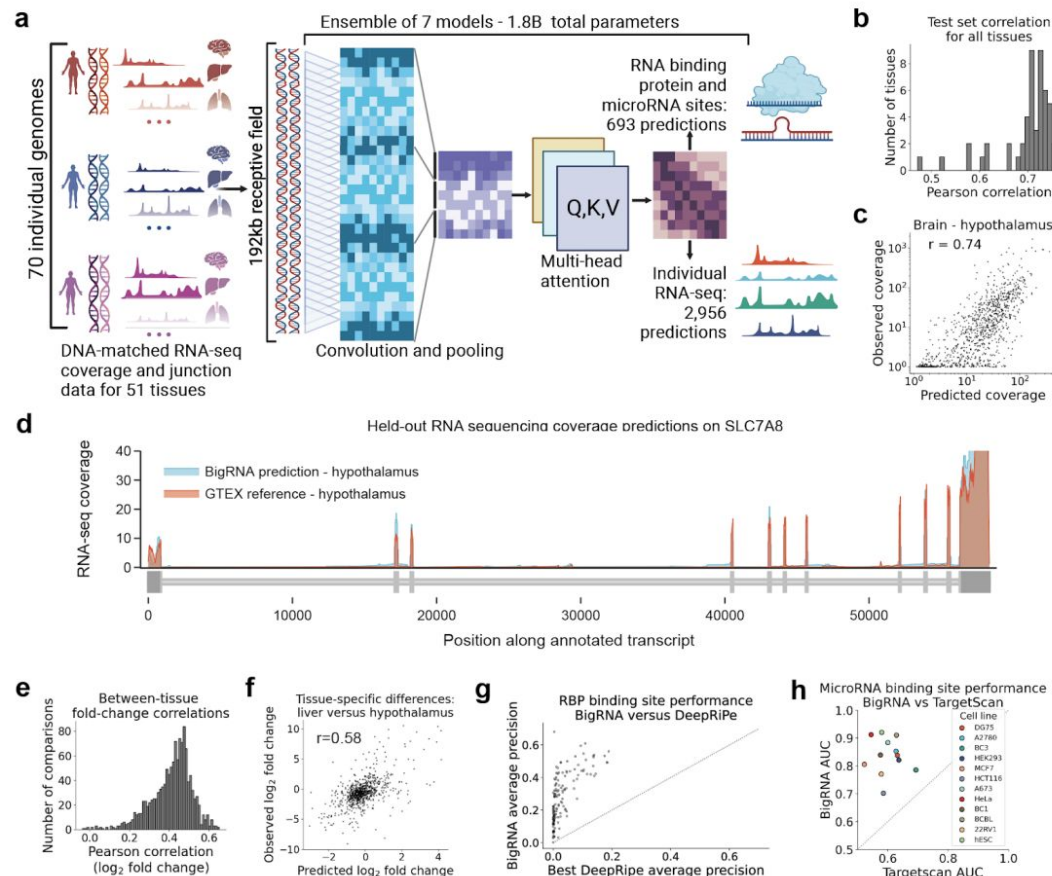


<https://github.com/calico/borzoï>



Modelling omics data (3)

BigRNA
GTEx



Thank you for listening

See you in 30 minutes for the practical session 💪🐦

Ensure you have a GitHub account, you will need it!

