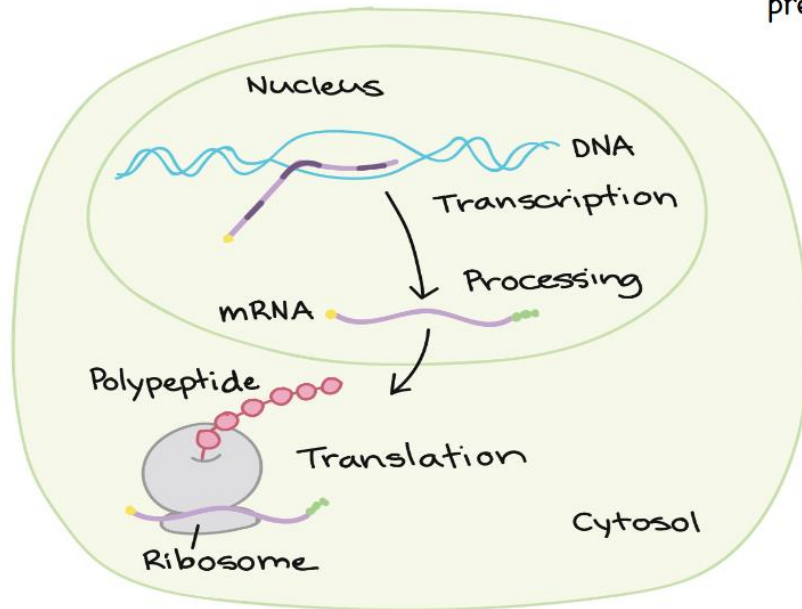


Background for RNA-Seq Data Analysis Workshop

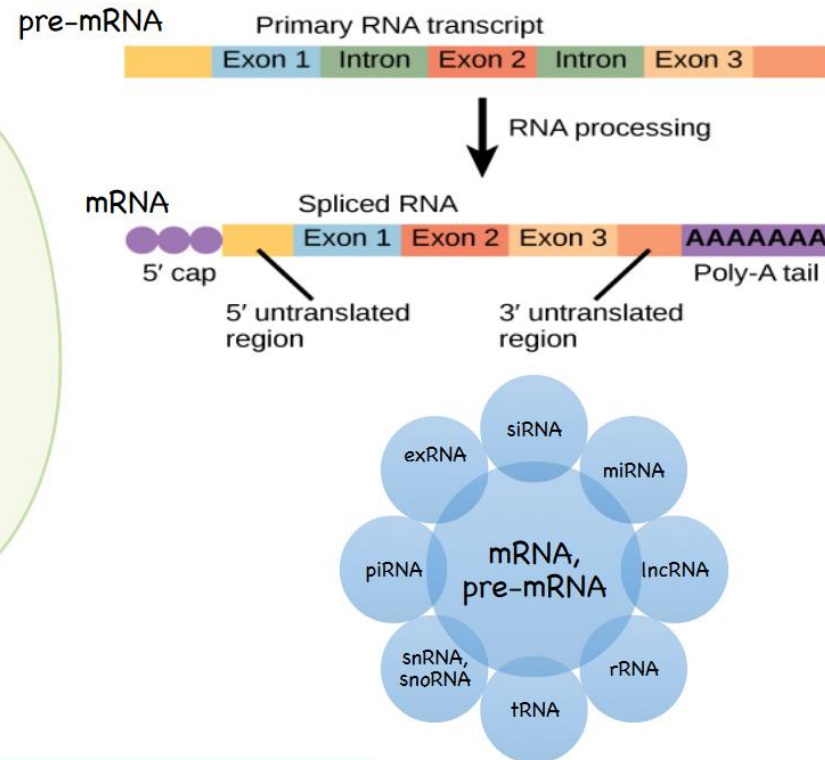
Bioinformatics Core

Wellcome Centre for Human Genetics

Profiling the transcriptome



- Image credit Khan Academy Open Courses



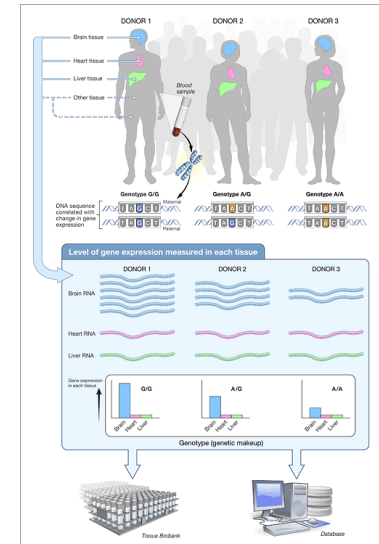
Gene expression profiling techniques provide valuable insight into complex biological systems, albeit a snapshot – highly dynamic and tightly regulated process (splicing, gene methylation, RNA stability/degradation, miRNA regulation etc.)

Wide utility of gene expression profiling

- ENCODE
- Allen brain atlas
- Genotype-Tissue Expression Project (GTEx)
- TCGA

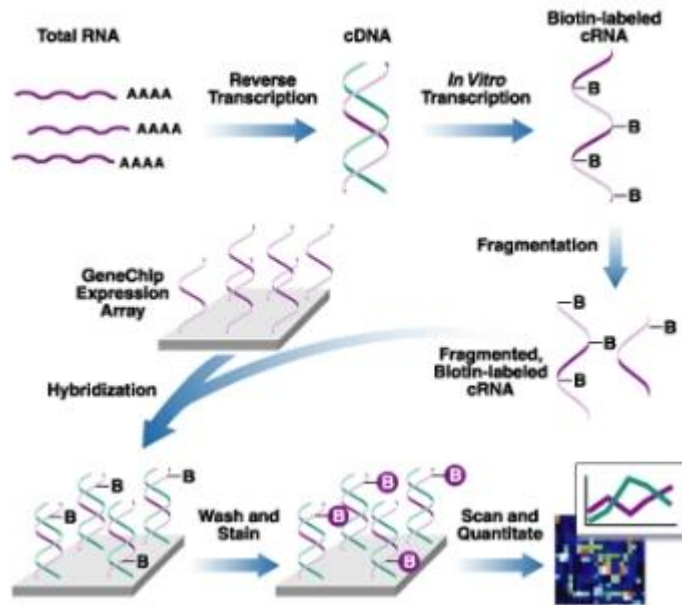


ALLEN BRAIN ATLAS
DATA PORTAL



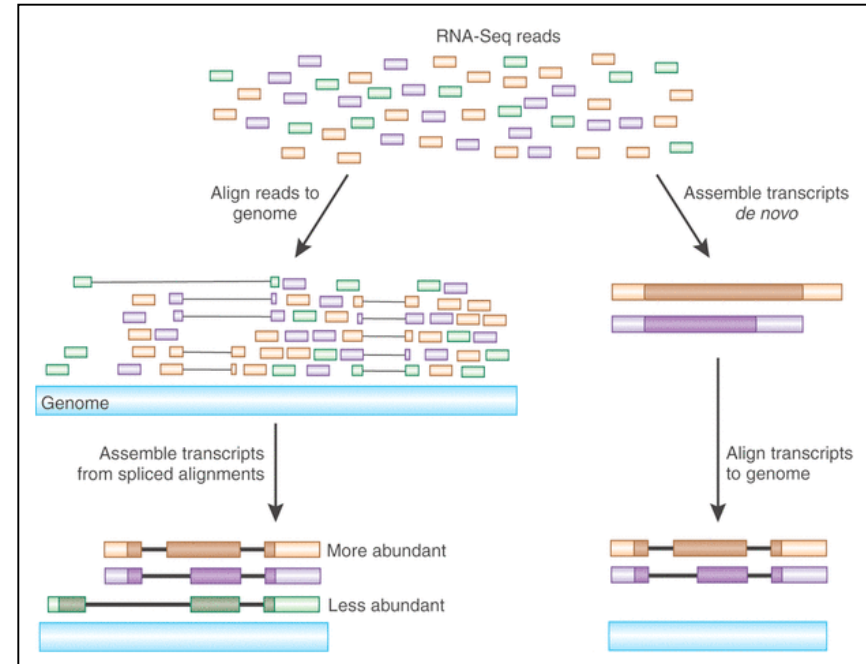
- Public repositories containing tens of thousands of datasets
 - Gene Expression Omnibus (GEO) <http://www.ncbi.nlm.nih.gov/geo/>
 - ArrayExpress <https://www.ebi.ac.uk/arrayexpress/>
 - Sequence Read Archive (SRA) <http://www.ncbi.nlm.nih.gov/sra>

Microarrays



Complementary hybridisation
early 1990s onwards

RNA-Seq



Next-generation sequencing
2007 onwards

- No prior knowledge of gene sequences needed in RNA-Seq; array probes have to be designed
- Better suited for transcript discovery, isoform characterisation and refining existing annotations e.g. length of 5' or 3' UTR, uncharacterised exons
- Costs quite comparable between technologies now (RNA-Seq used to be more expensive)
- Even though RNA-Seq produces count data, these do not represent individual transcripts but many fragments of transcripts
- Short-read platforms like Illumina are not capable of sequencing entire transcripts with a single read, so isoform reconstruction and splicing characterisation are very challenging tasks
- Myth: RNA-Seq can detect low expressed genes better than arrays
 - In most experiments, up to half of all genes are not sequenced at all or generate just a handful of reads
 - Additional sequencing will still tend to be from more highly expressed genes, so lower end extremely hard to interrogate
- Caveat: what you sequence in an RNA-Seq library influences your data for all genes – very inter-dependent in a way that arrays are not

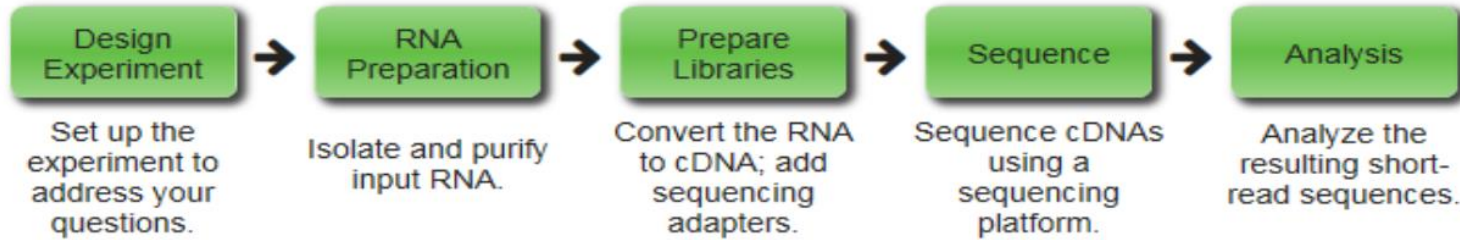
RNA-Seq: Library preparation

RNA-Seq library preparation protocol	Pros/Strengths	Cons/Weaknesses	Requirements
PolyA enrichment	Selects mature mRNAs and full transcript covered	Requires good quality samples	≥ 200 ng input
Ribo-depletion	Retains wider range of RNA species Full length of transcripts covered	Higher sequencing depth required, increasing cost	
3' mRNA	Very cost-effective; works well for variable quality/quantity input material	Can't characterise isoforms/splicing	
SMARTer low input	Good when only small amounts of RNA can be obtained	More expensive	10ng input
Small RNAs (miRNA)	Additional layer of informative data	Needs separate library prep so an extra cost	
Single cell	High resolution data, novel insights	Expensive, data analysis considerations	pg



- Number of reads required per sample depends on selected protocol and experimental questions
- HiSeq4000 – one lane = 250 million reads
- Multiplexing e.g. 10-plex human samples gives ~25m reads for each, plenty for quantifying gene expression (for those genes that are expressed)
- Higher depth required in some situations e.g. for splicing analysis, certain library prep methods (ribo-depletion)

RNA-Seq: Pipelines



- ✓ Multitude of algorithms and pipelines available.
- ✓ Most approaches correct, but have to be tailored to the needs of the investigators in order to better capture the desired effect.



Adapter trimming



Alignment



Assembly



Quantification



Count normalization



Alternative splicing/
Isoform level analysis



Differential
expression analysis

- Any sensible pipeline will produce reasonable results
- **DON'T** worry about trying to get to a definitive answer.....
- Instead, make sure appropriate tools are selected for the task and then carefully used:
 - Is the tool suitable for my question and data I have?
 - Have I understood how it works and how the parameter settings affect its behaviour?
 - Have I provided the right input and made sure the output is sensible?
 - Have I checked my R code for mistakes or unintended behaviour?

It is all too easy to for untoward things to happen somewhere along the line, and also surprisingly hard to spot them in high dimensional data like RNA-Seq

Raw gene counts are influenced by:

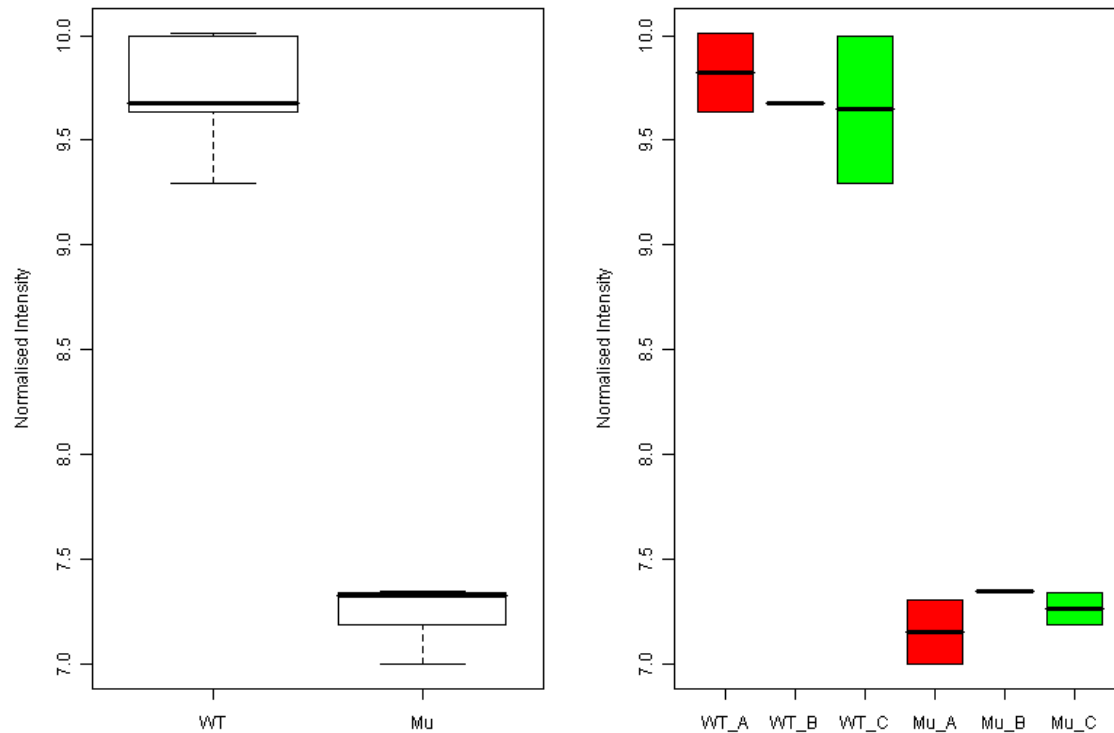
- Sequencing depth - varies by sample
- Gene length
- Amplification/sequencing biases – GC content
- Positional biases – preferential locations for RNA fragmentation
- Also by sample composition, which is less obvious

Count summaries accounting for various factors

- RPKM (reads per kilobase per million reads)
- FPKM (fragments per kilobase per million reads)
- CPM (counts per million reads)
- TPM (transcripts per million reads)

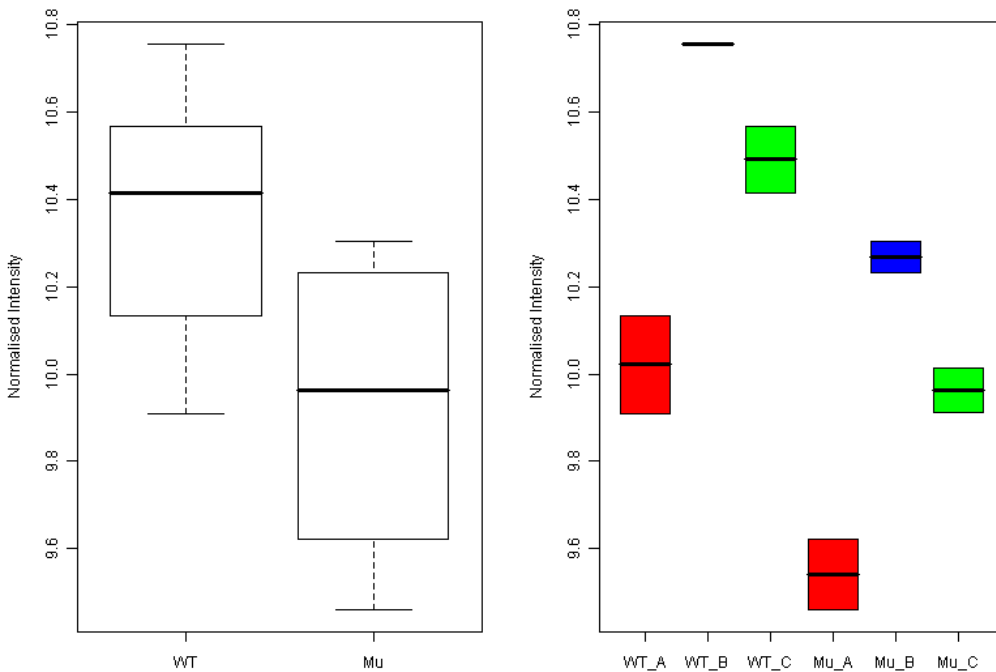
- Gene expression data highly sensitive to many factors
 - Lab operator/conditions, day performed, sample collection methods, RNA extraction day and so on
 - Often influence the data to greater extent than any experimental effects
 - Any step where treated and control samples are handled differently could confound the experiment
 - If split into batches containing mix of treated/control samples, can account for potential effects in analysis
- Also be aware of potential effects from factors unrelated to the experiment on the data, which may need to be accounted for to optimise analysis

Example differentially expressed gene



- Wt and Mut groups
- Three different litters (A, B, C)
- Top gene ~ 5x higher expression in Wt compared to Mut
- Similarly expressed across litters in both genotypes

Gene with strong litter effect



- Within litters, consistent pattern of higher expression in WT vs Mut
- Within genotypes, expression depends on litter: B > C > A
- Accounting for this source of variability increases power to detect changes of interest

- A typical RNA-Seq experiment has very small sample size, maybe 3-6 replicates per condition
- Small samples generally make it harder to infer differential expression, because there is more uncertainty in the estimation of the mean and variance, which are in turn used to compute the test statistic and p-value.
- Unreliable estimates of the sample variance due to small sample size can be problematic - potentially producing false positives when the variance is under-estimated and reducing power when the variance is over-estimated.
- Built into limma and edgeR packages is a procedure to 'borrow information between genes' and improve the robustness of the gene variance estimates
 - Variance estimates are moderated towards a common value according to how other genes with similar expression levels behave

- Whichever technology used, inherently limited to descriptive results no matter how well experiment performed or data analysed.
- Produce large amounts of information; subjective interpretation, can be mined in different ways requiring human decision-making to take the results further.
- What genes/pathways to focus follow-up experiments on? Different researchers could easily identify quite different themes from the same results.
- Best used as starting point for further work - following up hypotheses from gene expression data to uncover mechanistic/causal effects can produce elegant studies.