

Small RNA transcriptomics

IN-BIOS5000/9000

Trine B Rounge



Learning outcomes

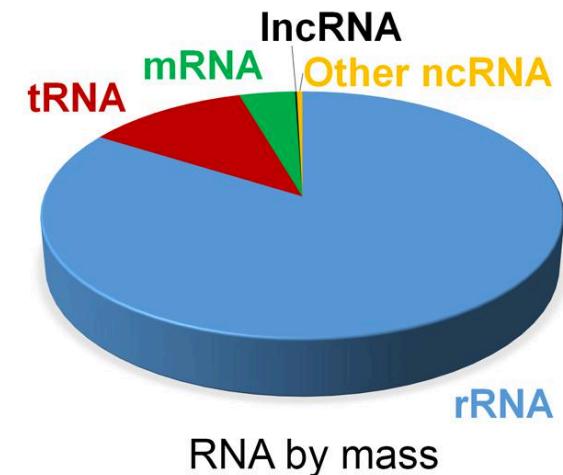
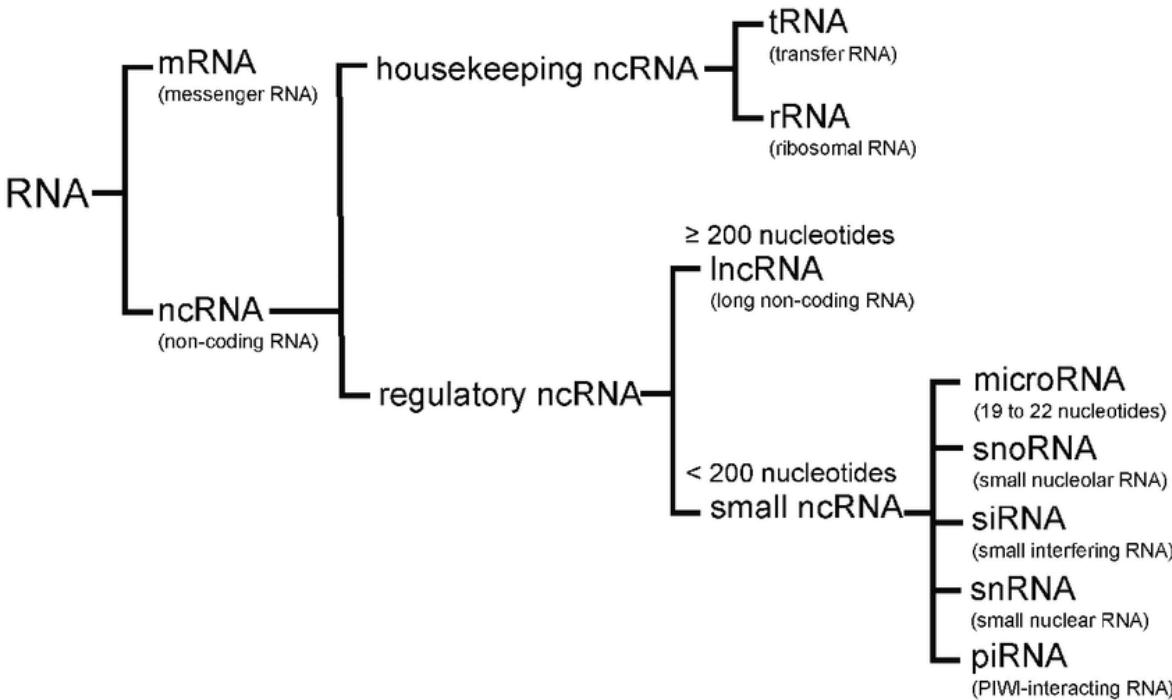
1. What are small RNAs
 - Their role in the cell and in bodily fluids
2. What is small RNAs transcriptomics
 - Methods/technologies
 - Experimental design
3. Analyse small RNA-sequencing data
 - Read Counts
 - How to identify differences
4. Research examples (if time)
5. Practical tasks



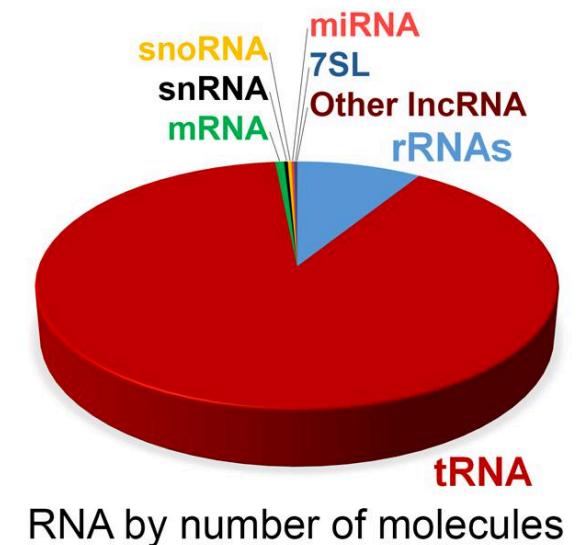
1. Small RNAs



What are small RNAs



RNA by mass

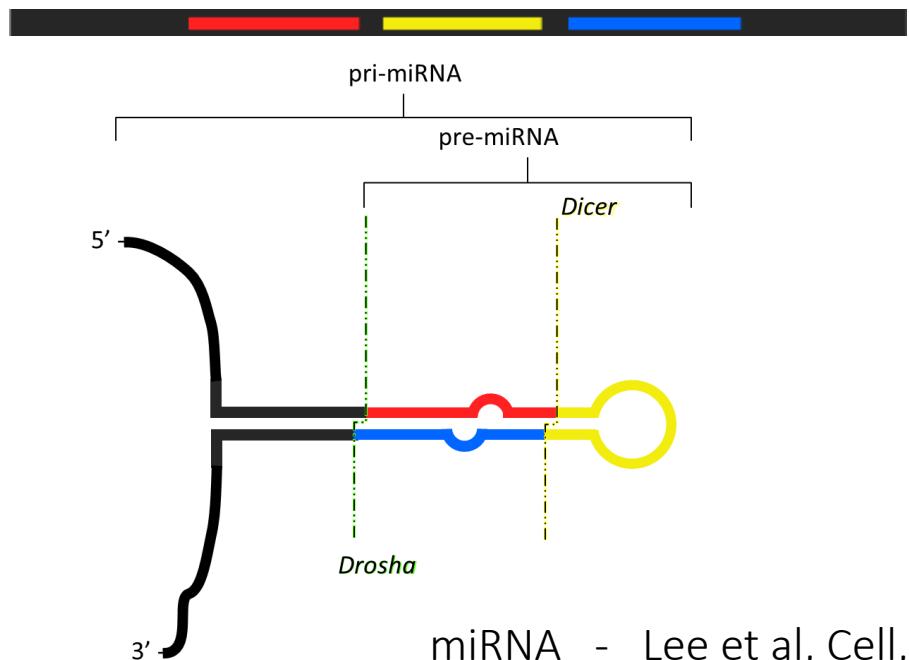


RNA by number of molecules

Palazzo & Lee, Front Genet, 2015
Small RNA transcriptomics

miRNAs

~ 22 nucleotide in length

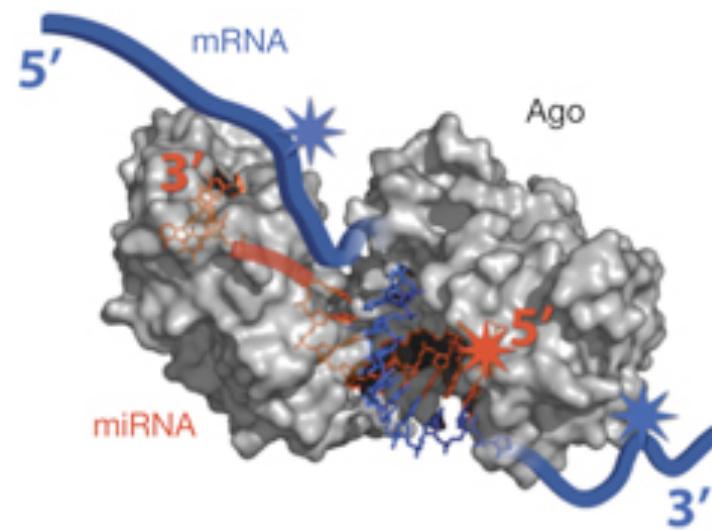


miRNA - Lee et al, Cell, 1993

Dicer - Hutvagner et al, Science, 2001

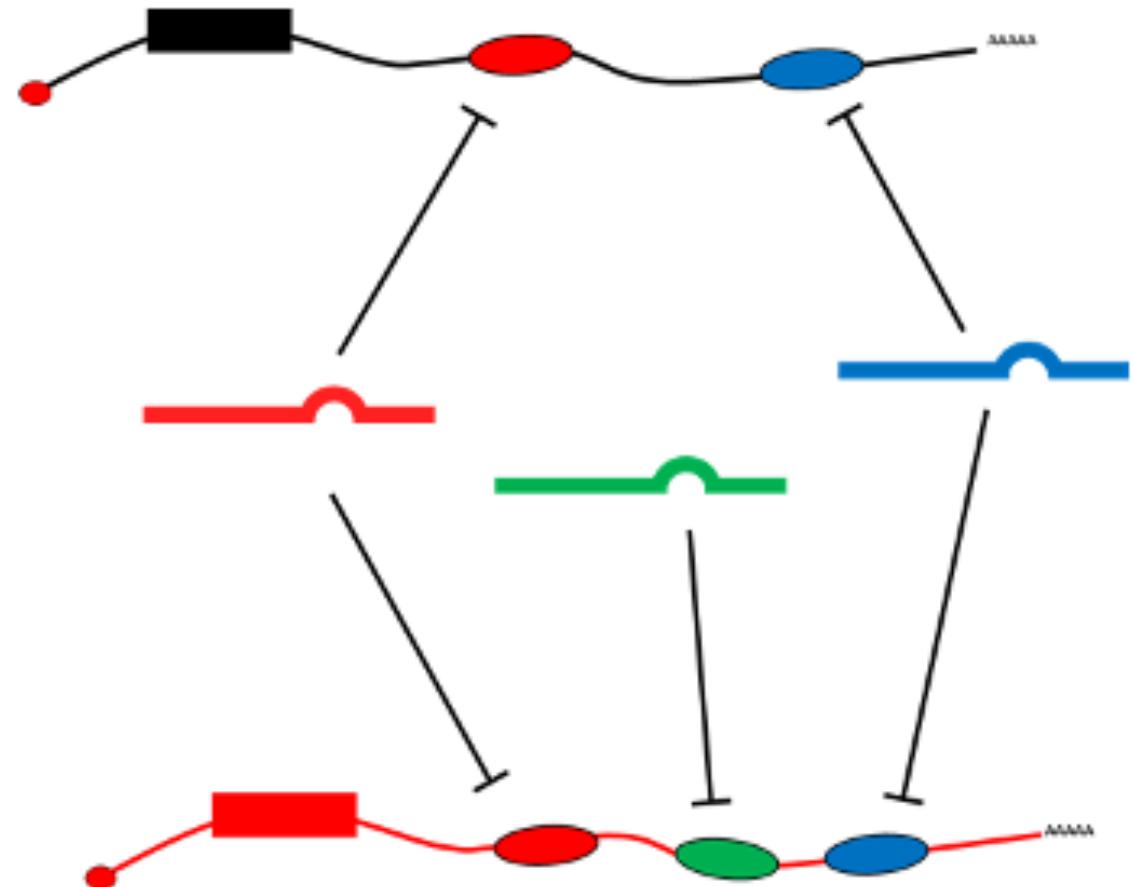
Drosha - Lee et al, Nature, 2003

RNA-induced silencing complex (RISC)
including the Argonaute protein



miRNA function

- The miRNA-RISC complex targets mRNA for silencing.
- target multiple mRNAs and multiple miRNAs can target the same transcript.
- the fine tuning of most protein products within the cell.



Pichler & Calin, Br J Cancer 2015

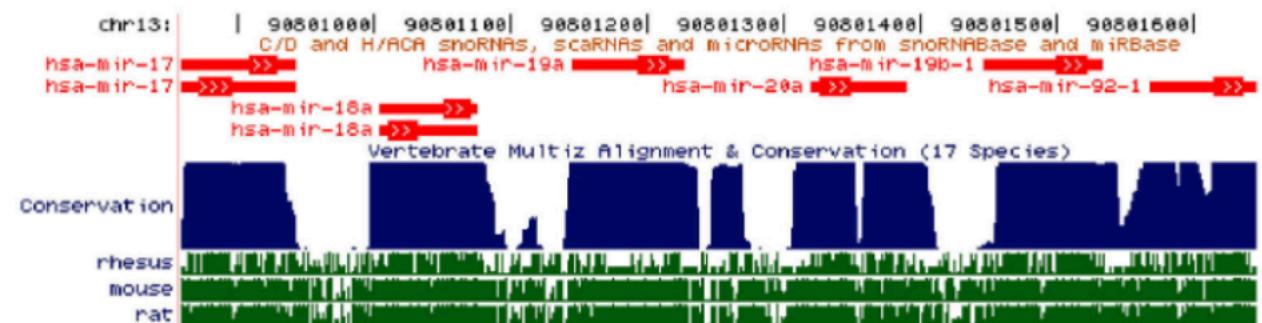
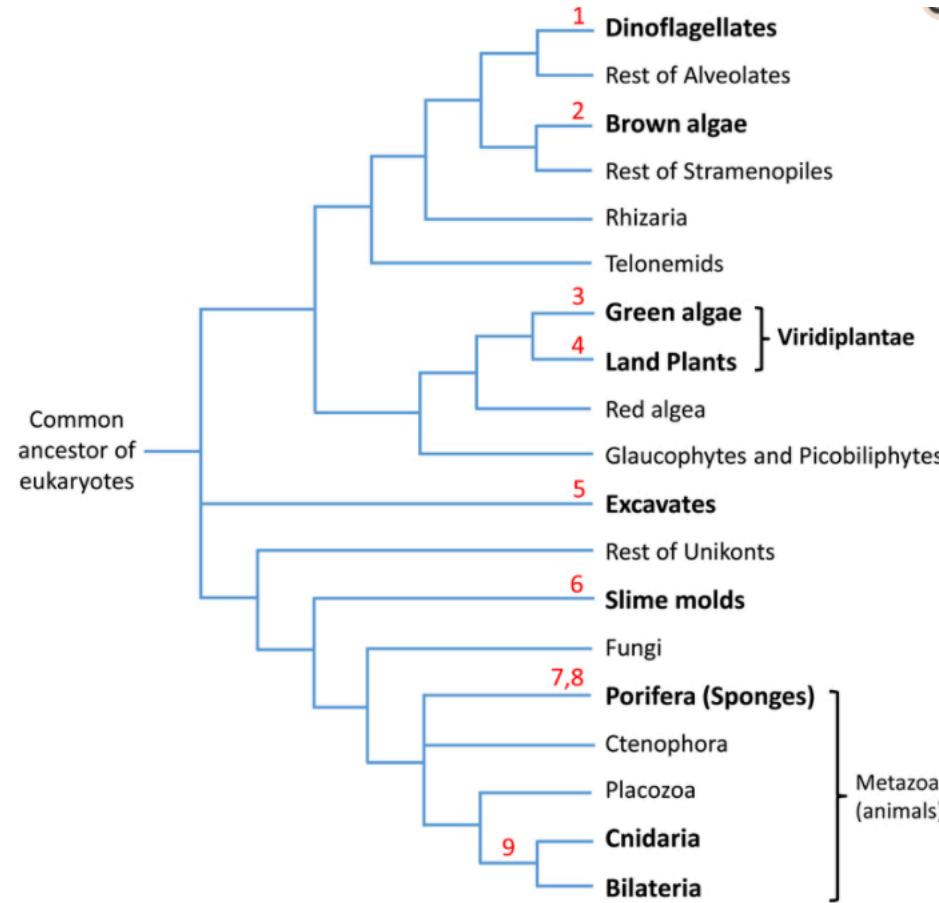


UiO • University of Oslo

IN-BIOS5000/9000

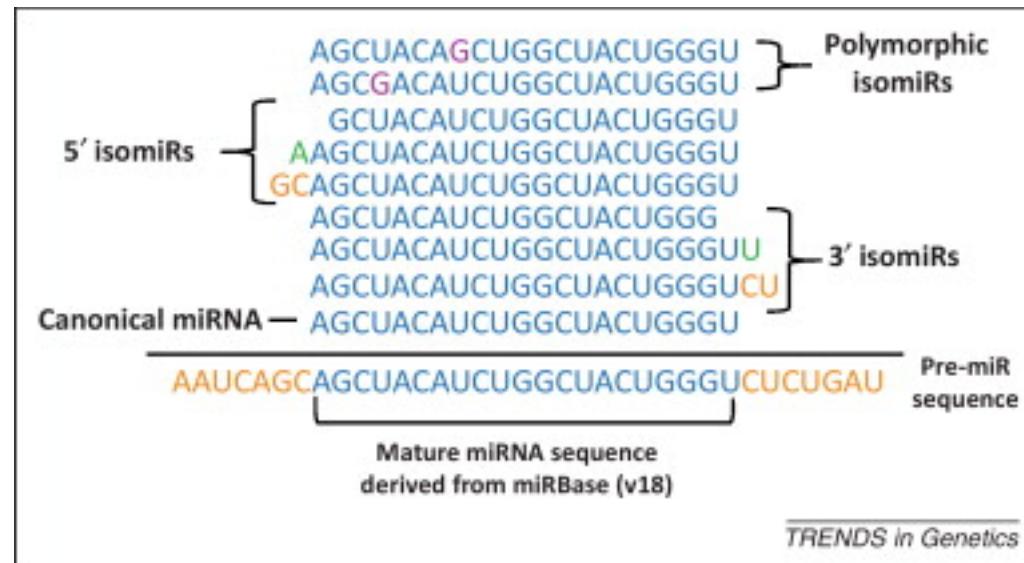
Small RNA transcriptomics

Evolutionary conserved



Isoforms

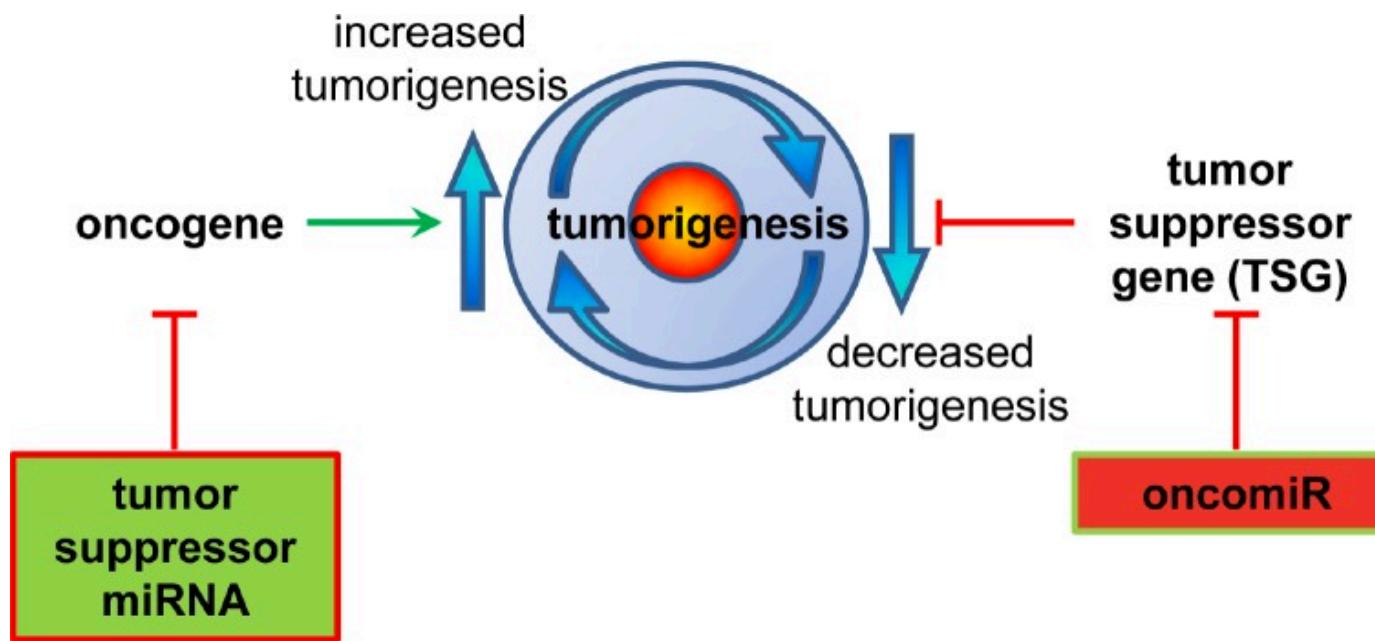
IsomiRs



Variability in Dicer and Drosha processing

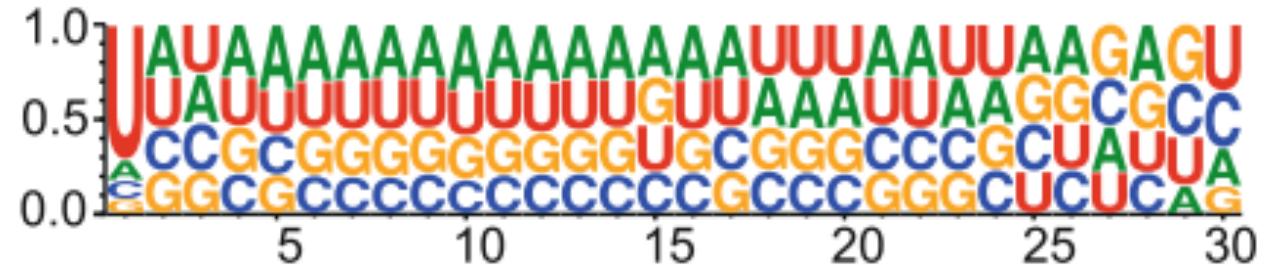
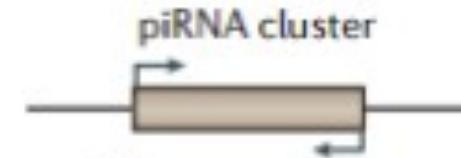


miRNA expression in cancer



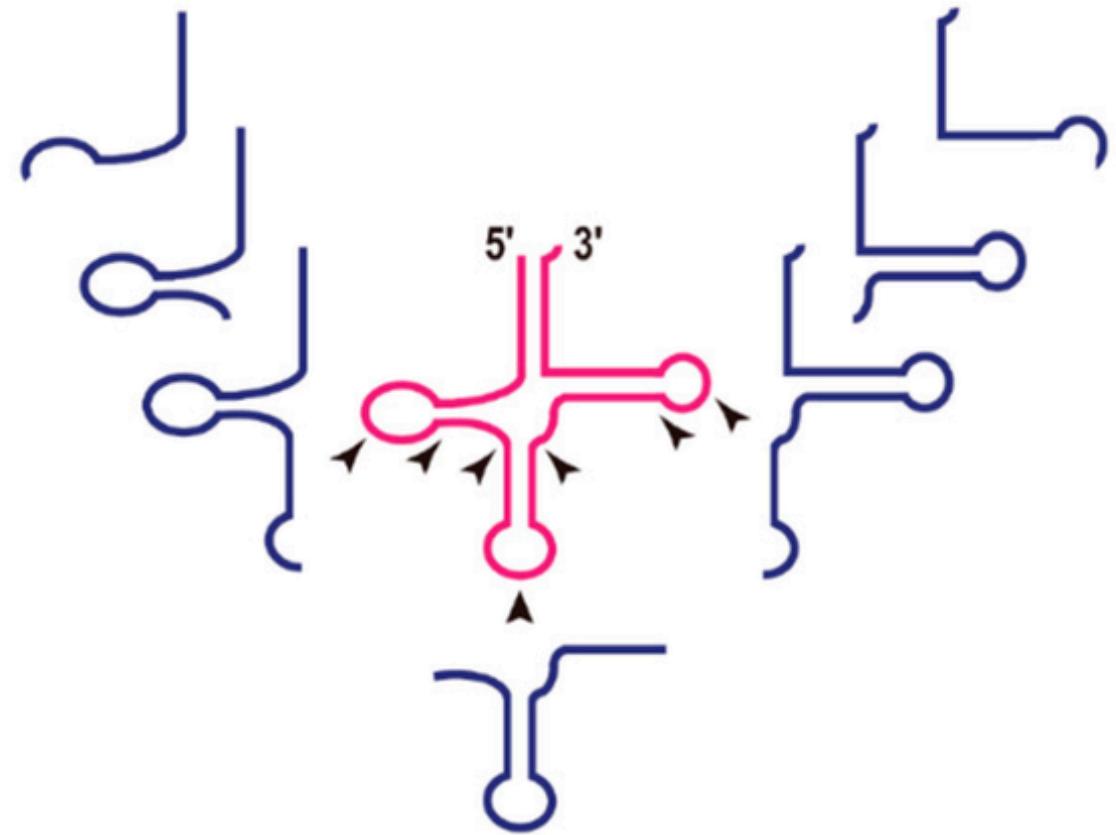
piRNA

- 26-31 nucleotides in length
- transcribed in clusters
- 5' uridine
- a role in RNA silencing via PIWI
- active in the testes of mammals
- silencing of transposons
- silencing via RISC
- amplified by a ping-ping mechanism



tRNA fragments

- 76 - 90 nucleotides in length
- carrying an amino acid to the ribosome
- Fragmented in:
 - halves
 - stress induced
 - fragments (1/4)
 - RISC regulation,
 - epigenetic control,
 - metabolism,
 - immune activity
 - stem cell fate commitment



snRNA and snoRNA

- snRNA – small nuclear RNA
 - Role in splicing
 - ~ 150 nucleotides.
 - Bohnsack et al, Biol Chem, 2018
- snoRNA – small nucleolar RNAs
 - role in the modification, maturation, and stabilization of rRNA
 - regulation of gene expression and alternative splicing
 - stress response
 - Liang, Frontiers in Oncol, 2019



2. Small RNAs transcriptomics



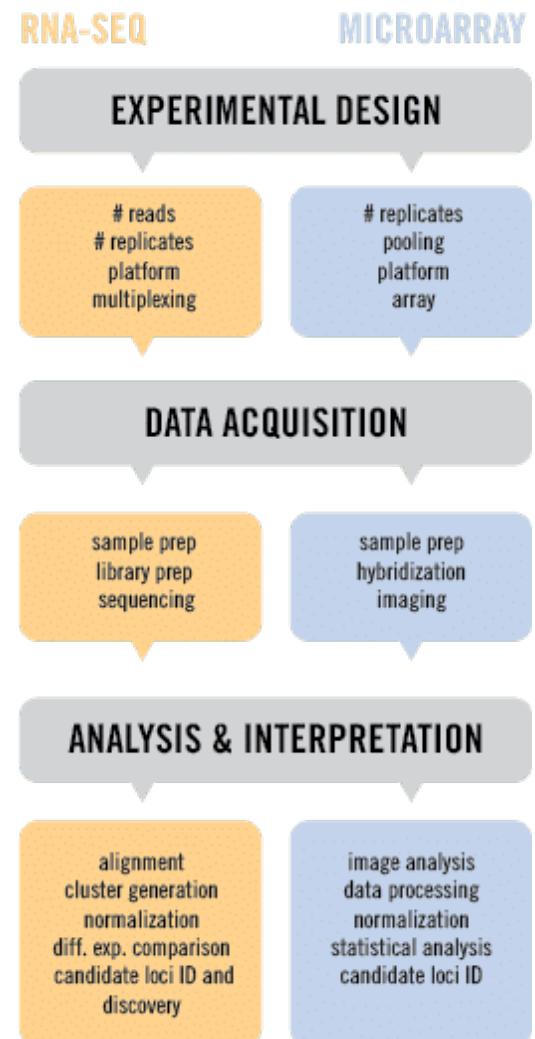
Small RNA transcriptomics

- Study of the the **complete set of small RNA transcripts** that are produced by the genome, under specific circumstances or in a specific cell or body fluids—using high-throughput methods
- <100 nucleotides or <200 nucleotides
 - Transcribed as small RNAs
 - Processed to small RNA
 - Degraded to small RNA

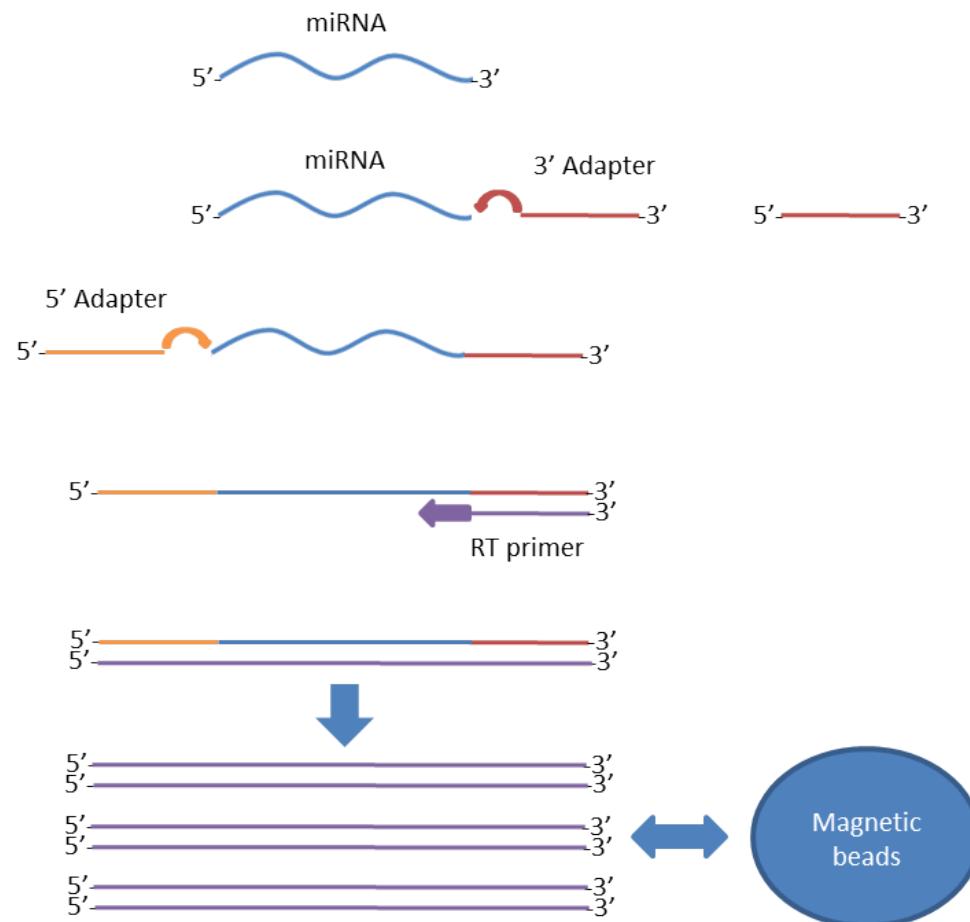
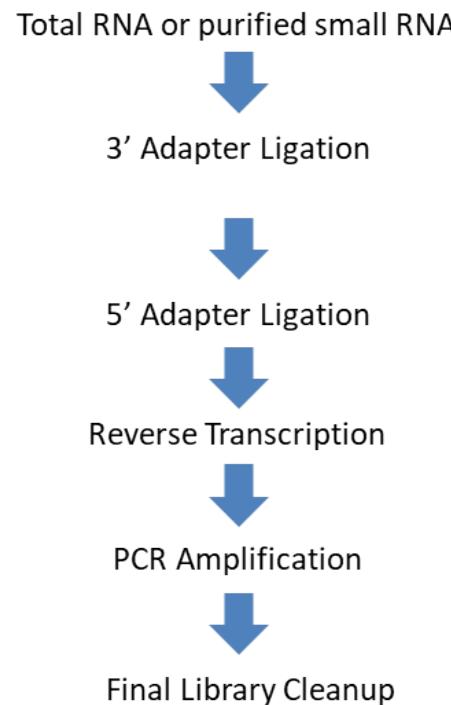


Small RNA Technologies

- RT-qPCR panels
 - limited part of the transcriptome
 - good reproducibility
 - good for validation
- Microarray
 - limited to known genomes and known transcripts
 - In some cases poor sensitivity and a limited dynamic range
 - non-specific hybridization or cross-hybridization
 - well defined RNAs
 - low cost
 - low noise
 - MirNome
- Small RNA sequencing
 - potential ambiguity in the mapping
 - protocol specific biases
 - discovery
 - large range
 - Many kits – not comparable results

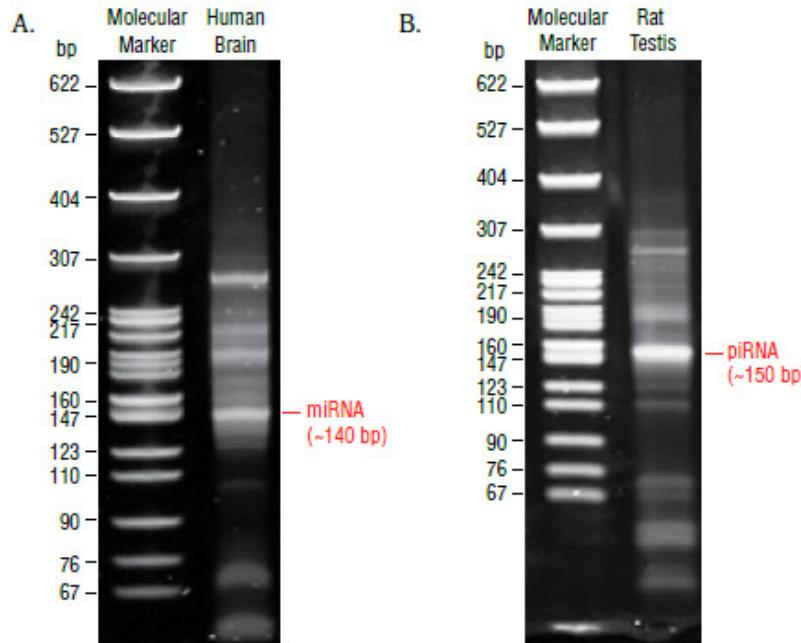


Small RNA sequencing – library preparation

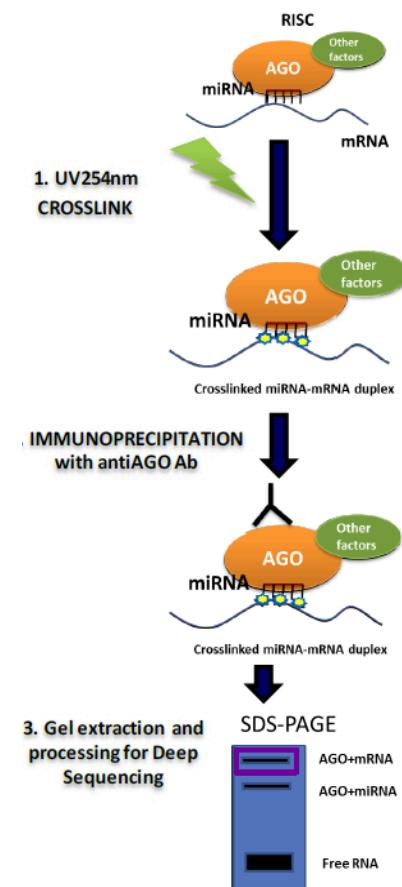


Small RNA transcriptomics

Size selection

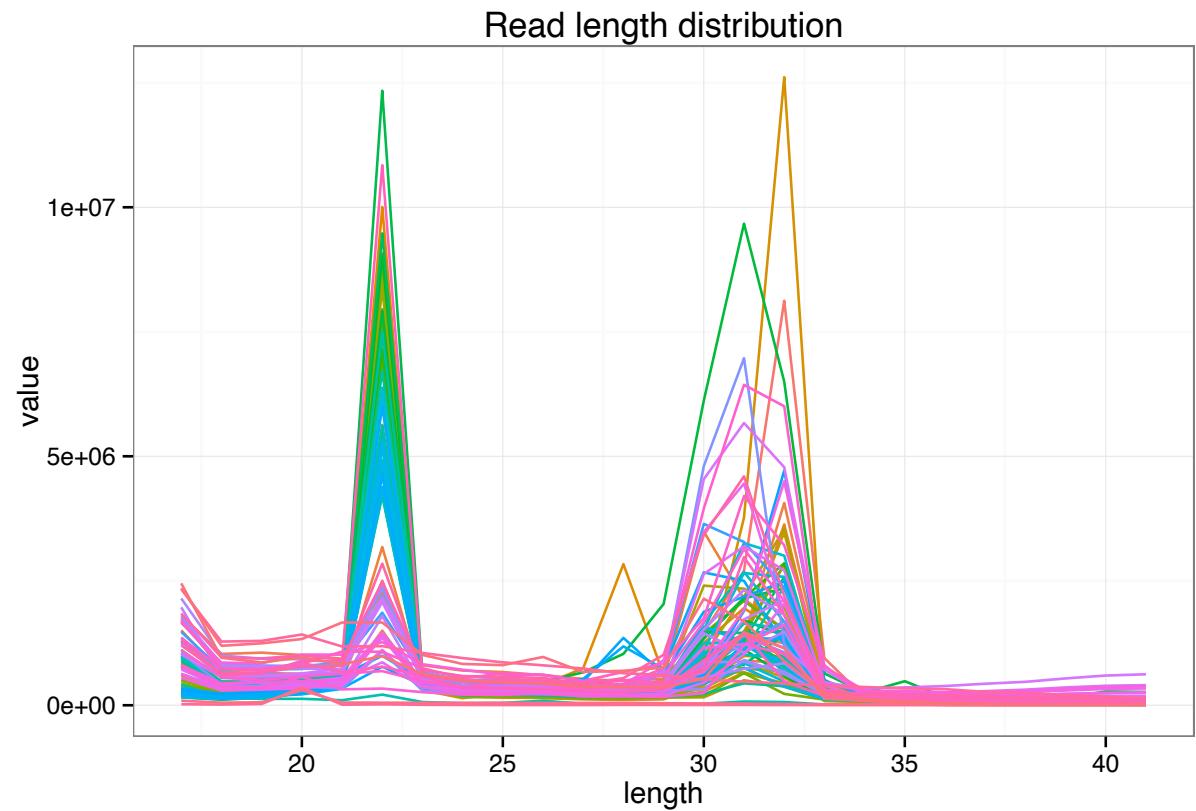


AGO immunoprecipitation



Small RNA transcriptomics

- Size selection is **not** perfect
- Bias can be introduced
- Size distribution need to be checked



Sequencing

- The most common technology for small RNA-seq
- Reconstruction of isoforms
- Detection of novel transcripts
- Suited for **expression analysis**
- High number of samples



Sample size vs depth

Depth ↓

- RNA sequencing
 - 1. Highly expressed known transcripts
 - 2. Novel isoforms
 - 3. Low expressed / rare transcript

Sample size →

To detect a 1.5 logfold difference:

- 3 samples/group → 43% statistical power
- 10 samples/group → 91% statistical power

With 10 sample/group:

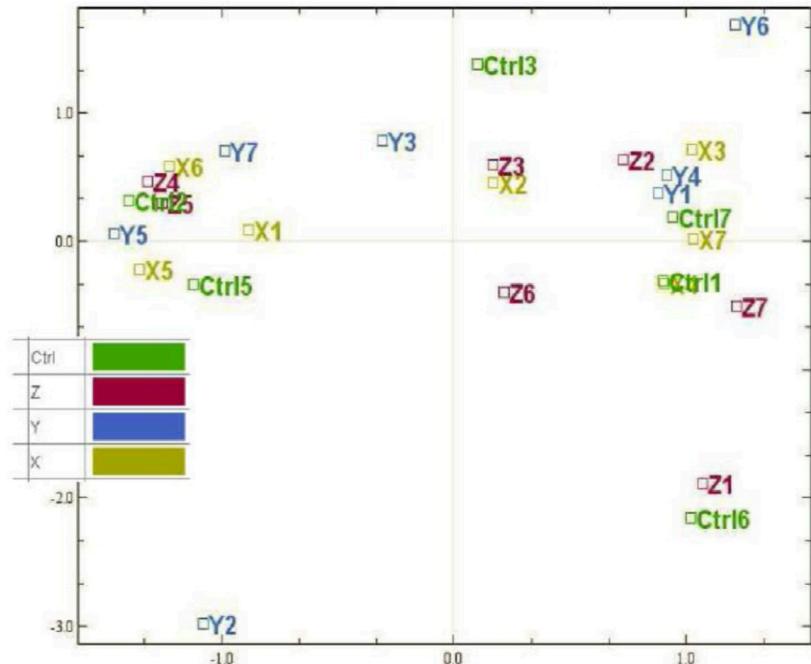
- 3 mill read/sample → 52% statistical power
- 10 mill read/sample → 80% statistical power

Talk to a statistician

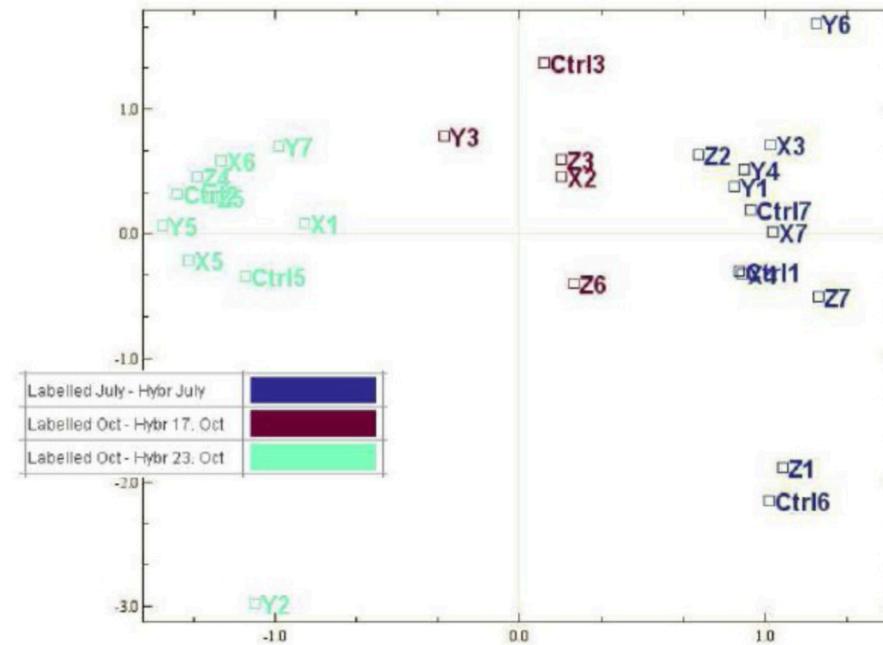
A survey of best practices for RNA-seq data analysis
<https://doi.org/10.1186/s13059-016-0881-8>



Batch effects



Samples color coded according to biology



Samples color coded according to labeling date

Mentimeter - quiz



3. Analyse small RNA sequencing data

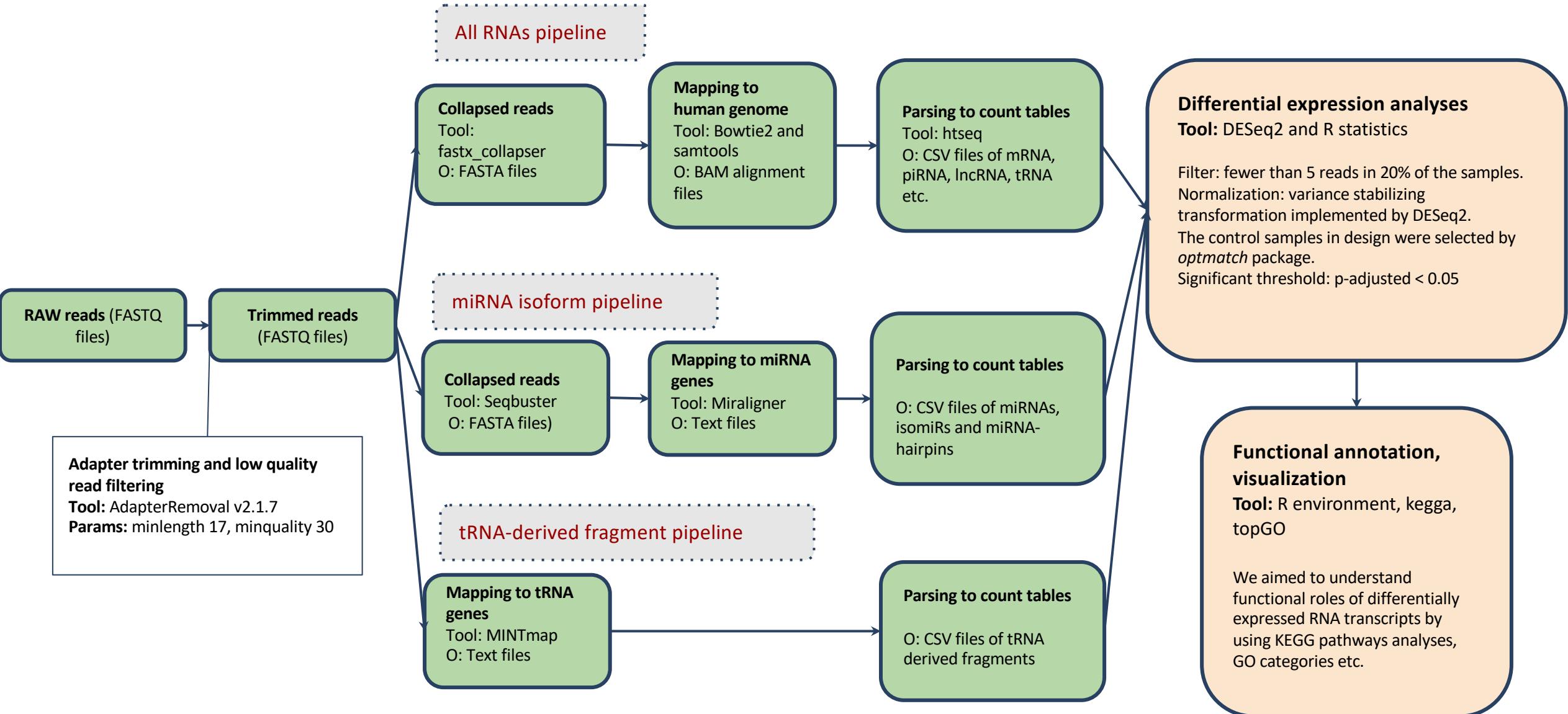


Small RNAseq workflow

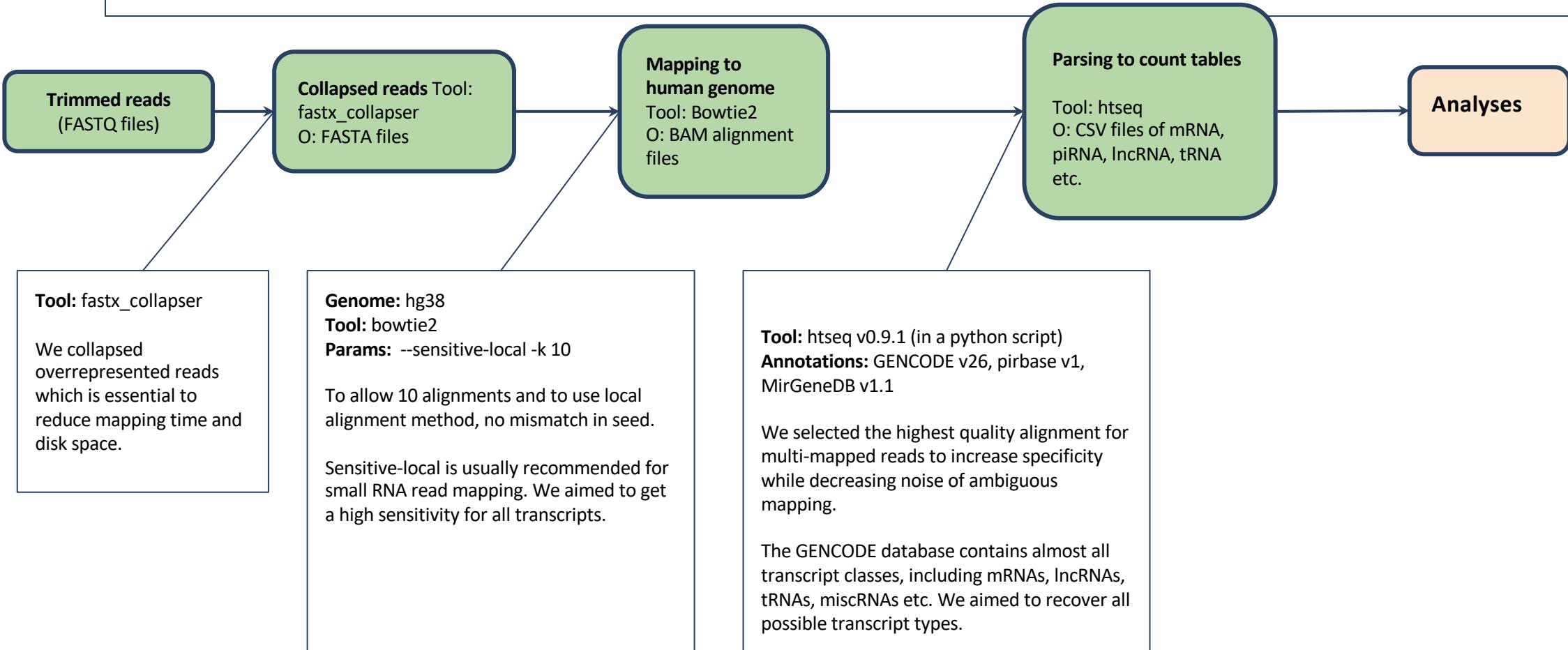
From Fastq files to small RNA read counts

- Reproducible
- Detect isoforms
- Detect most small RNA classes
- Efficient

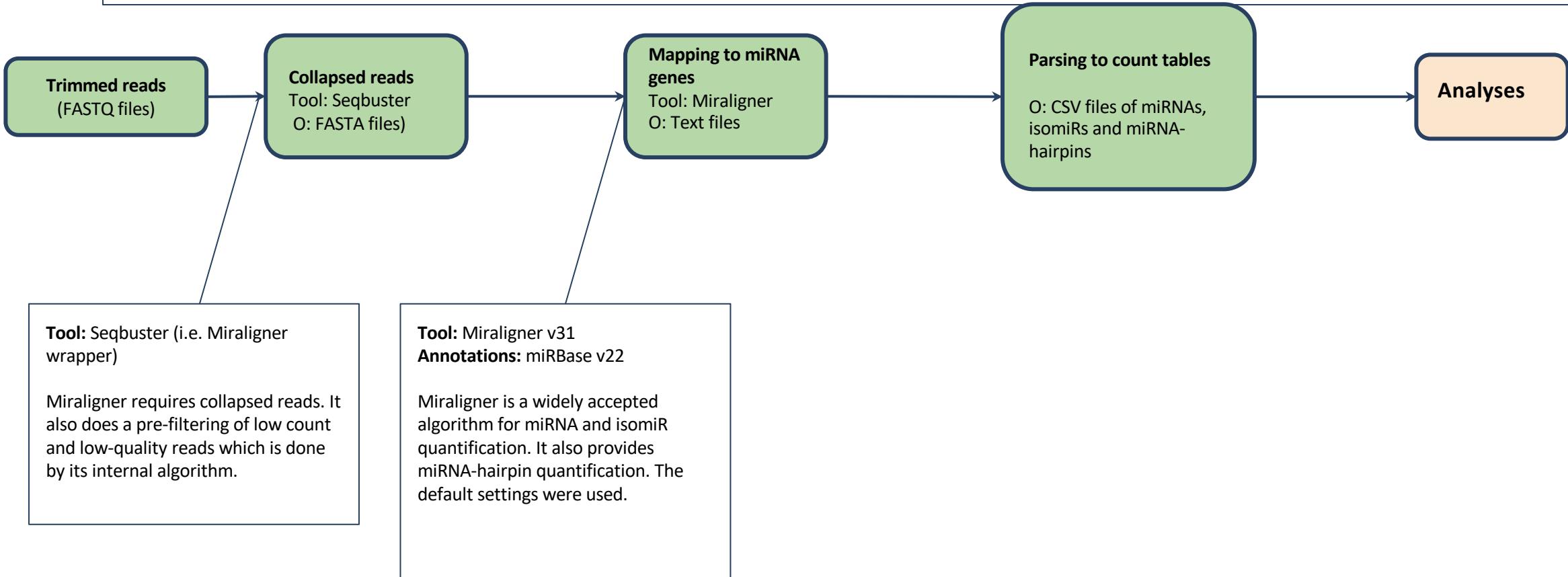




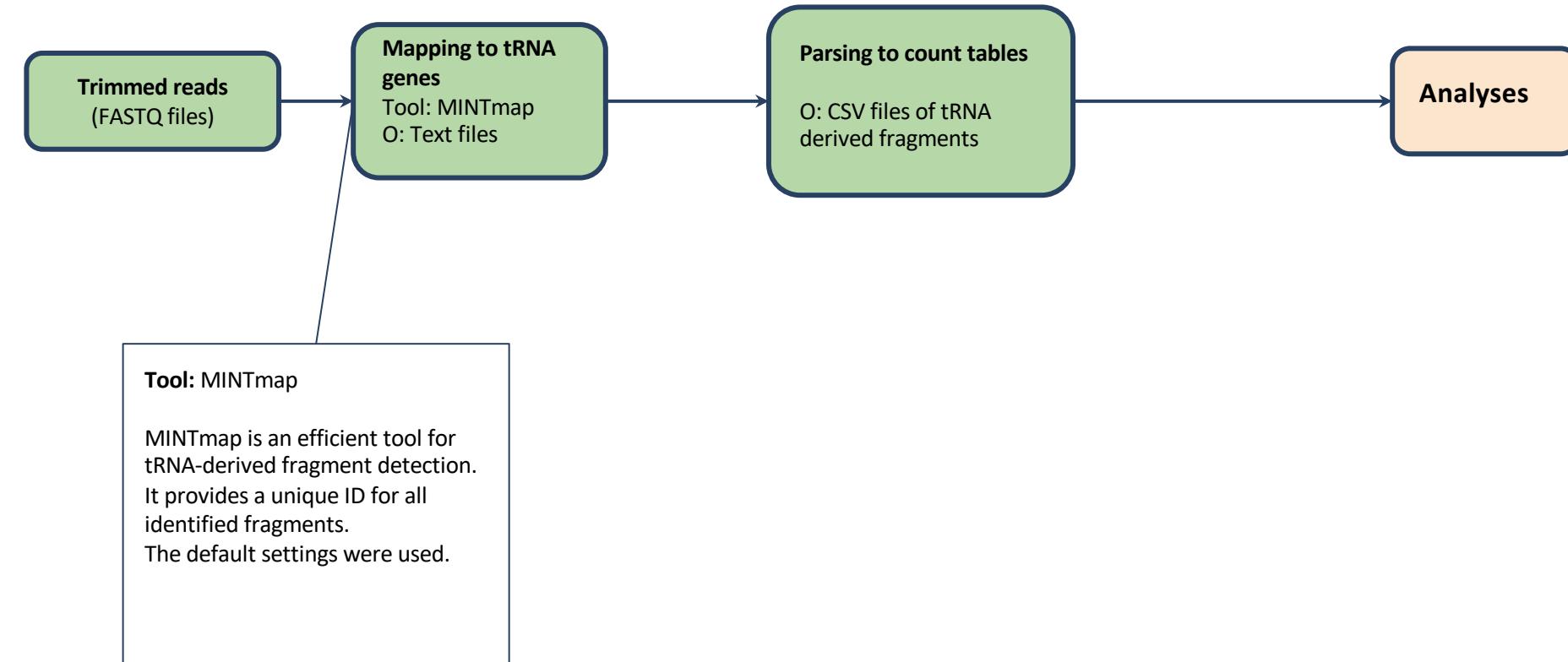
All RNAs pipeline



miRNA isoform pipeline

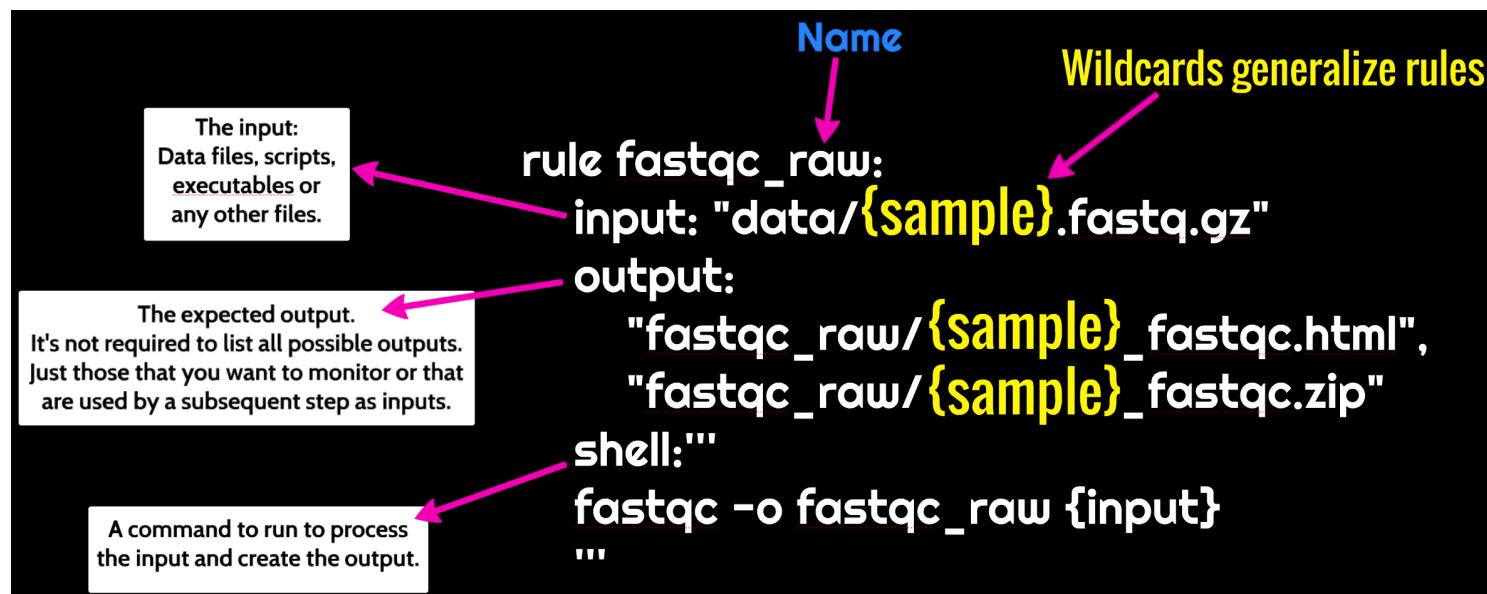


tRNA-derived fragment pipeline



Snakemake

- tool to create **reproducible** and **scalable** data analyses.
- python based language.
- Can be scaled to server, cluster, grid and cloud environments.



<https://snakemake.readthedocs.io/en/stable/>

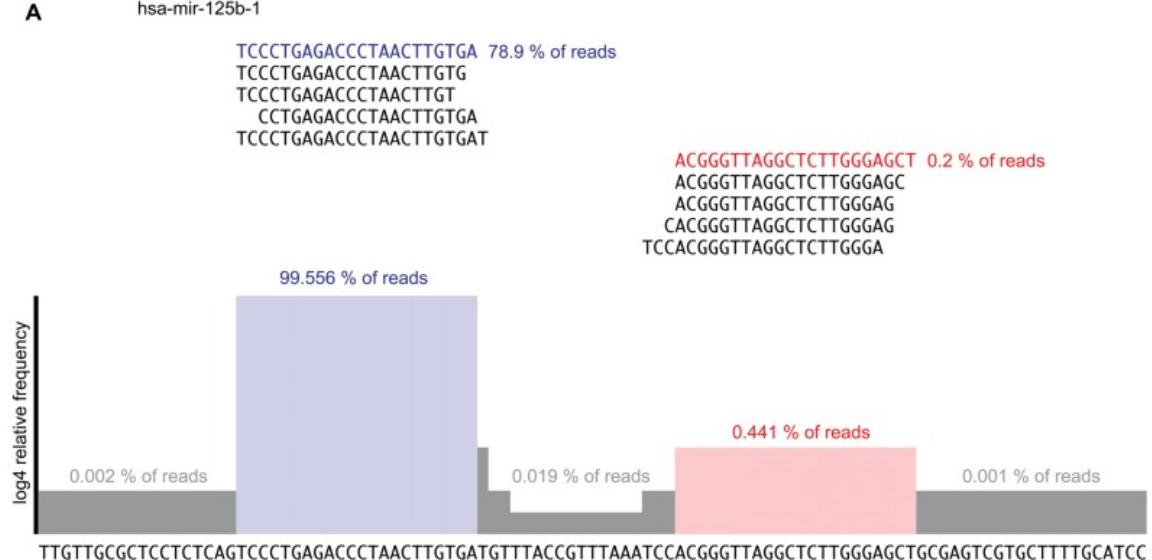
Small RNA count data

- Similar to mRNA count data
- Each class analysed separately
 - Detection/sequencing bias due to length differences

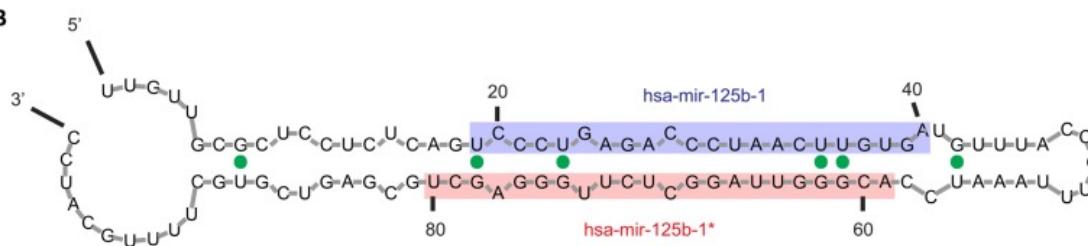


Read counts

A



B



Gene/Sample	Sample 1	Sample 2	Sample 3 repA	Sample 3 repB
Gene A	5	0	45	101
Gene B	17	500	32	67
Gene C	752	16432	20020	45078
Total	350250	278090	400890	799009

Normalization

IN-BIOS5000/9000

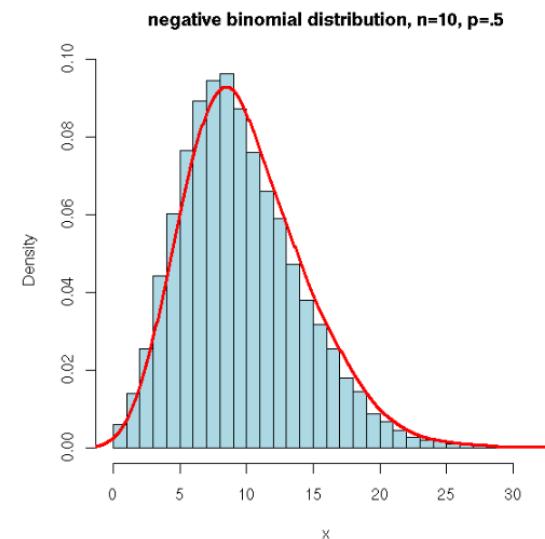
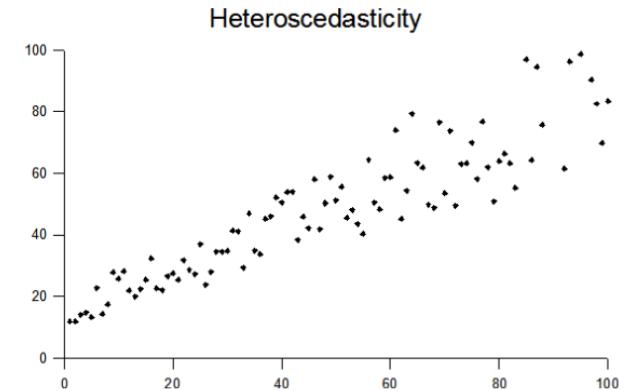


UiO : University of Oslo

Small RNA transcriptomics

Challenges of using count data

- The unit of measurement is number of reads (count) – not light intensities (microarray), or the shape of a curve (qRT-PCR).
- Large **dynamic range** – from zero up to millions.
- **Heteroscedastic** – variance is not equal across the range of counts
- Positive integers and non-symmetric distribution – i.e. normal distribution doesn't fit.
- Systematic biases (**normalization**)



Normalization

- Sequencing depth – library size
- Sequencing length
- Variance across means



Normalization for library size

- If sample A has been sampled deeper than sample B, we expect counts to be higher.
- Naive approach: Divide by the total number of reads per sample
- Problem: Genes that are strongly and differentially expressed may distort the ratio of total reads.



Normalization for library size

- To compare more than two samples:
- Form a “virtual reference sample” by taking, for each gene, the geometric mean of counts over all samples
- **DESeq2**: Normalize each sample to this reference, to get one scaling factor (“size factor”) per sample.

Anders and Huber, 2010
similar approach: Robinson and Oshlack 2010



Differential gene expression analysis

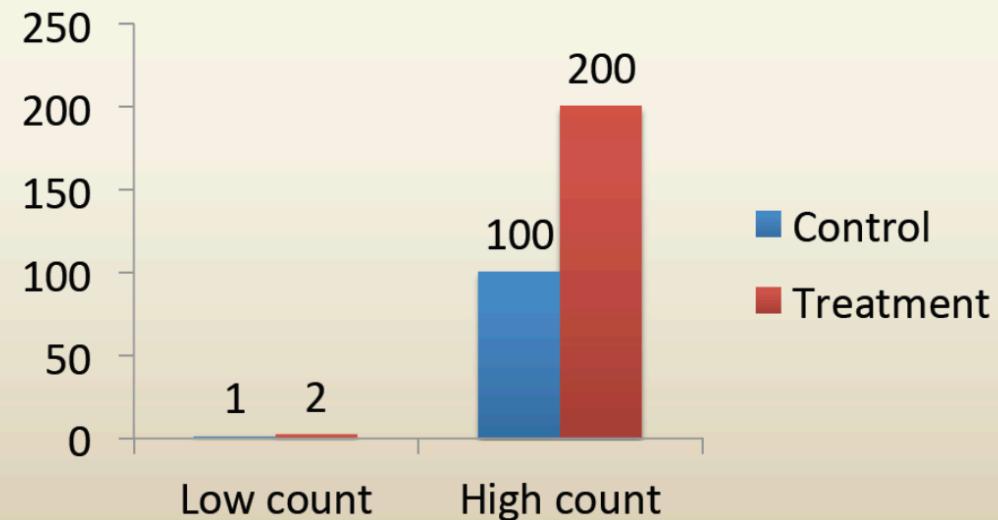
- DESeq2 implements a **Wald test**
- Conceptually similar to a t-test
- Goal is to identify genes that are differentially abundant (DE) between two conditions.
- Assumption: most genes are not DE
- Null hypothesis: each gene has the same abundance across conditions.



Strong poisson noise for low count values

I) Poisson counting error

- Uncertainty in count-based measurements
- Disproportionately large for low-count data



DE testing– adjusted p-values

Multiple hypothesis testing

- Thousands of genes = thousands of hypothesis tests (simultaneously)
- Increased chance of false positives! (Type I error)
 - e.g. you test for differential expression in 1000 genes that are not differentially expressed
 - You would expect $1000 \times 0.05 = 50$ of them to have a P -value < 0.05
- Individual P -values not useful
 - Need multiple testing statistic instead



Differential Expression

- **DESeq2** and **edgeR** – two of the most common packages for RNA-seq analysis (differential expression).
- **DESeq2** and **edgeR** – based on “raw counts” such as from HTSeq



Mentimeter - quiz



4. Research examples



Prediagnostic serum RNA dynamics in lung cancer



Phase 1

All LC samples
(542 cases and 519 controls)

Phase 2 - 7 time windows

All LC samples
(7 time windows, 27 cases and 135 controls each)

Phase 3 - 7 time windows

Early Stage
(9 cases and 45 controls)

Locally Advanced
(14 cases and 70 controls)

Advanced
(18 cases and 90 controls)

NSCLC
(84 cases)

NSCLC
(99 cases)

NSCLC
(165 cases)

SCLC
(75 cases)

ADC
(89 cases)

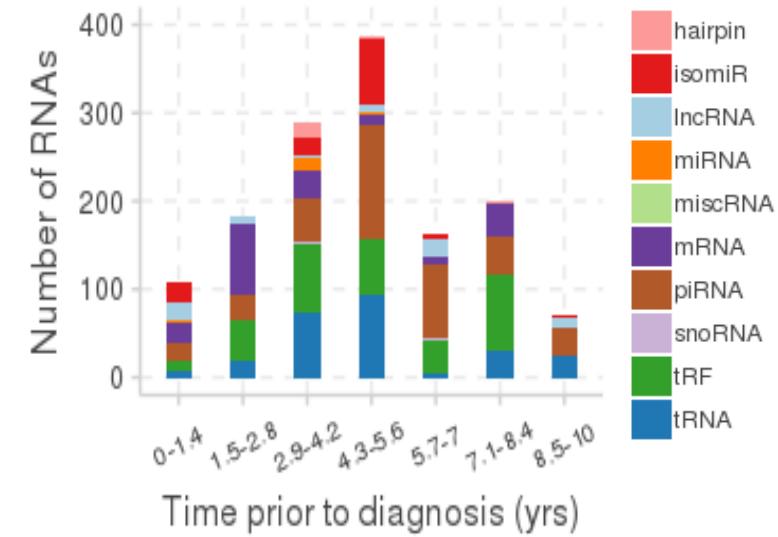
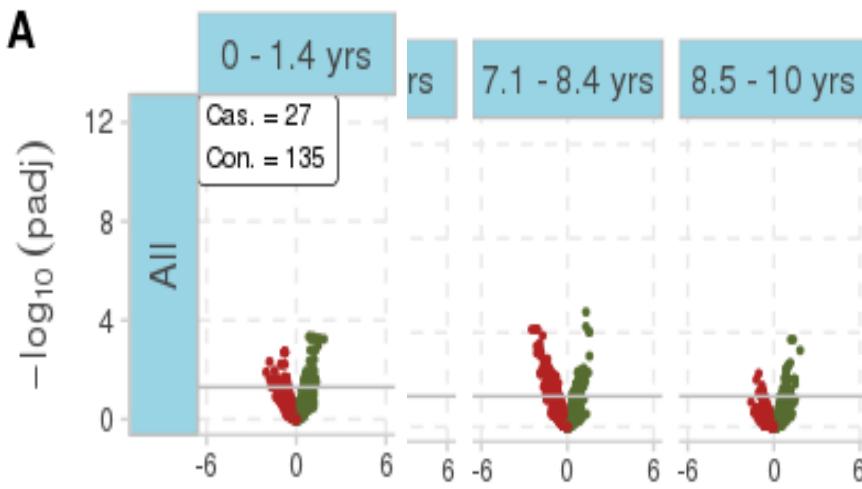
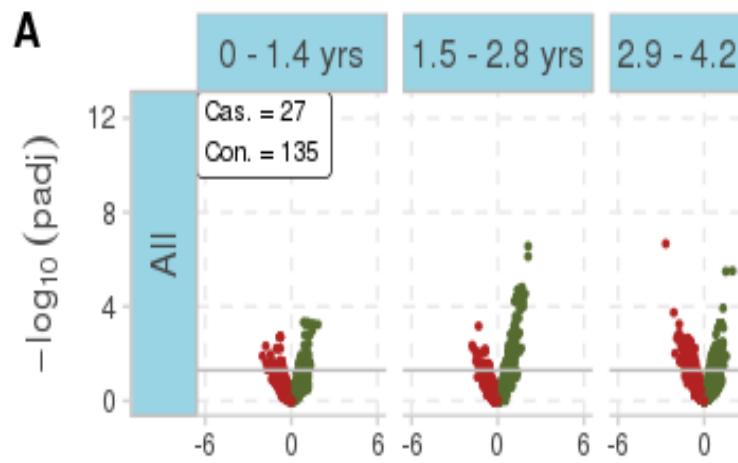
Phase 4 - sliding windows



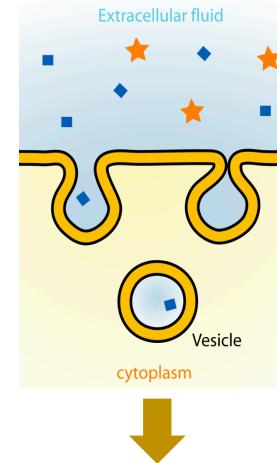
Sample and patient characteristics

			Stage			
	Control	Early (Localized)	Locally Advanced (Regional)	Advanced (Distant)	Unknown	
Histology						
NSCLC	-	84	101	171	11	
SCLC	-	9	35	76	4	
Others	-	10	5	32	4	

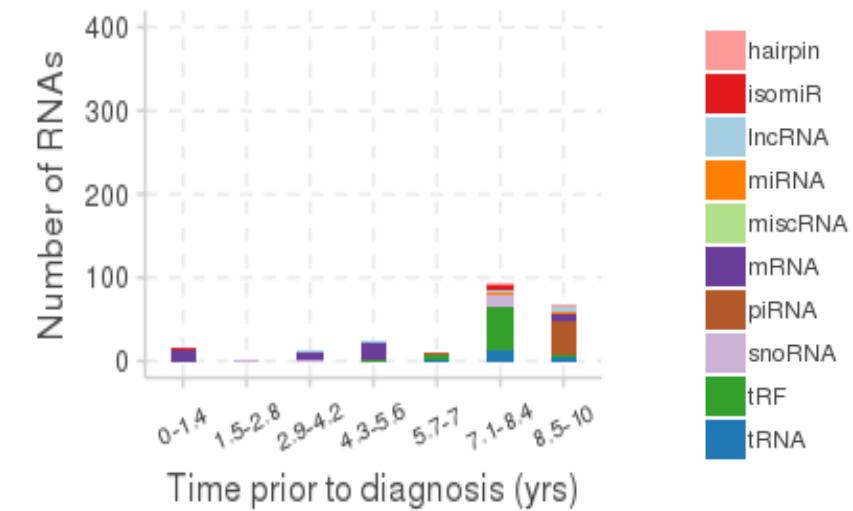
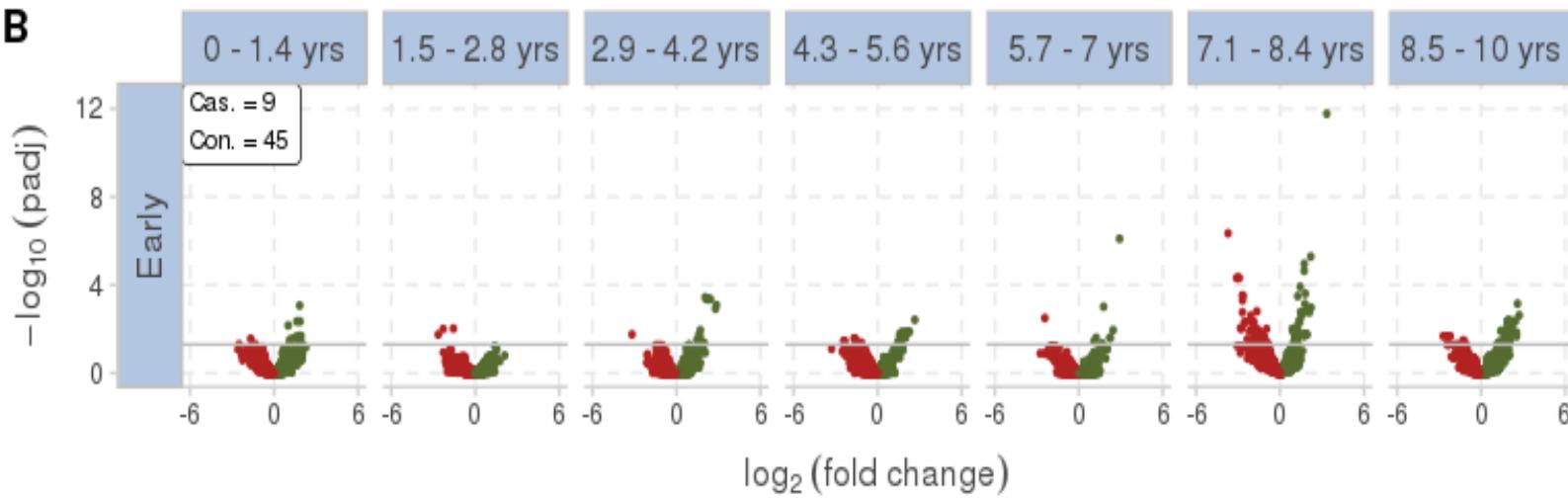




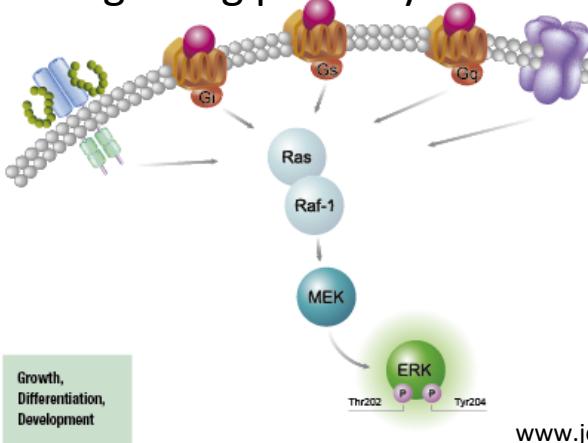
Endocytosis pathways



B



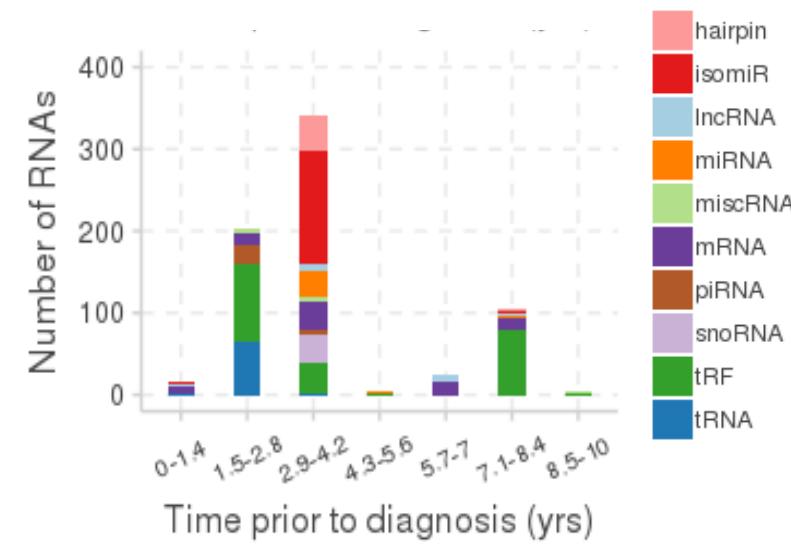
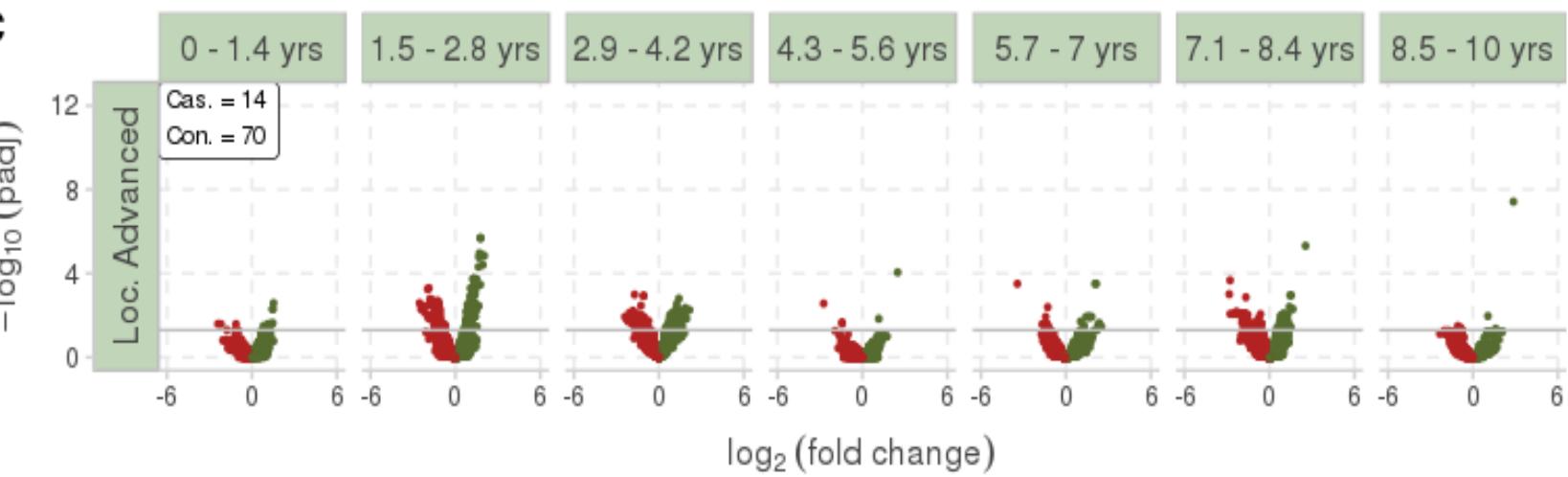
Signaling pathways



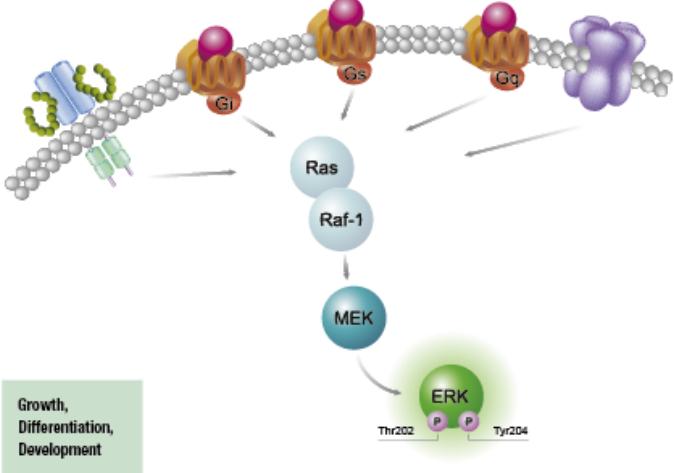
www.jcbio.co.kr/

MAPK, mTOR, ErbB, Ras

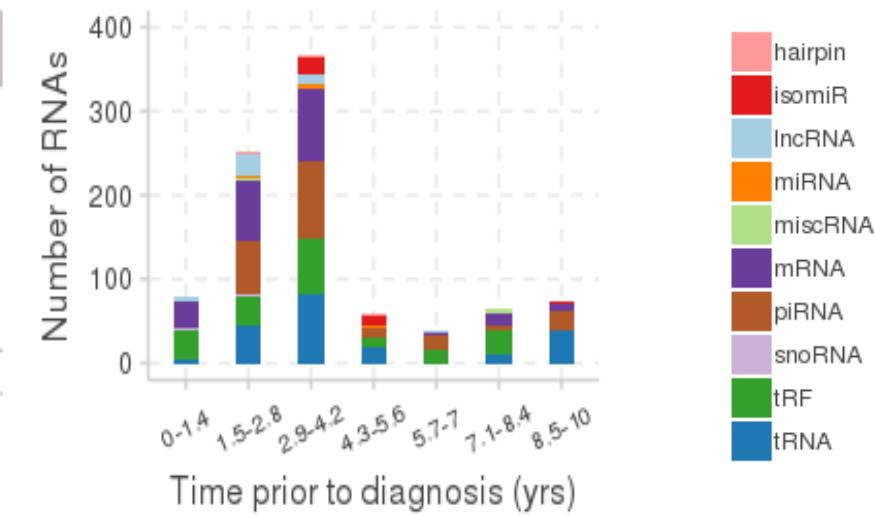
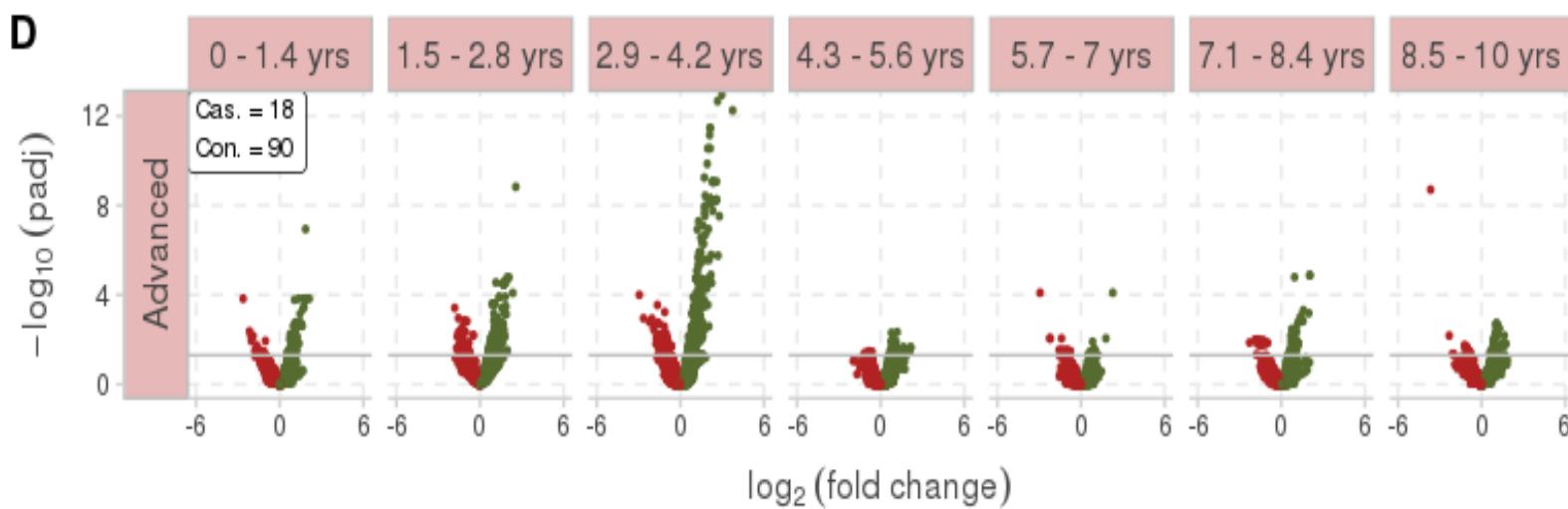
C



Signaling pathways



MAPK, ErbB,
Ras, mTOR
↓
MAPK, ErbB,
Ras
↓
ErbB
↓



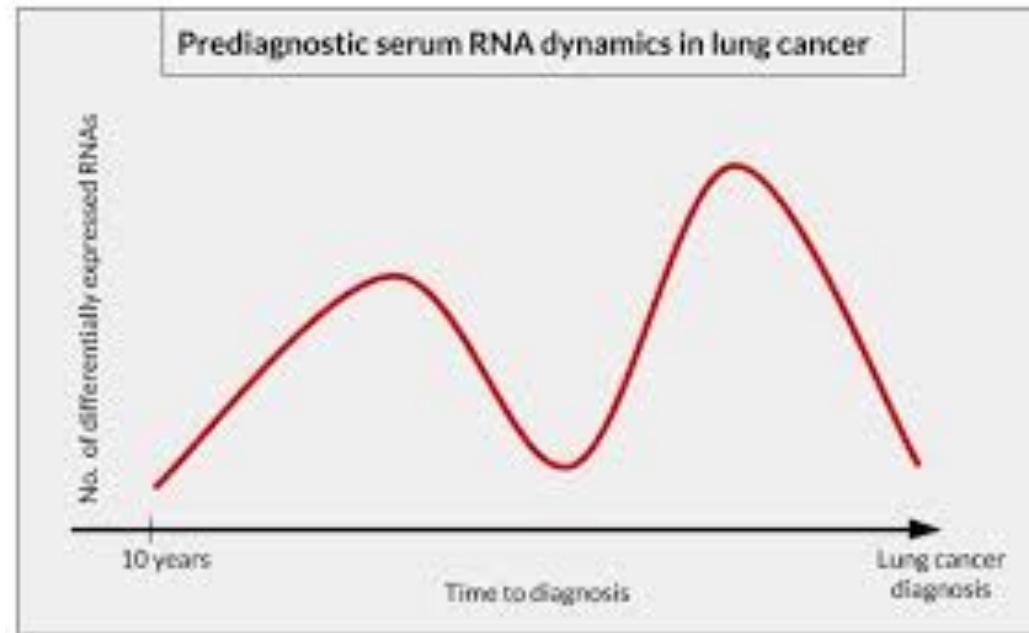
Summary

- Large-scale serum ncRNA analyses for biomarker research require
 - optimized methods
 - knowledge of technical and biological variation
- ncRNA signals in pre-diagnostic lung cancer serum samples are highly dynamic
 - Signals appear up to 10 years prior to diagnosis
 - Stage and histology specific
 - Disrupts proliferation related signalling pathways



A 10-year prediagnostic follow-up study shows that serum RNA signals are highly dynamic in lung carcinogenesis

Sinan Uğur Umu¹, Hilde Langseth¹, Andreas Keller^{2,3}, Eckart Meese⁴, Åslaug Helland^{5,6,7}, Robert Lyle^{8,9} and Trine B. Rounge^{1,10} 



Resources

- <https://edu.t-bio.info/course/transcriptomics-1/>
- <https://edu.t-bio.info/course/transcriptomics-2/>
- <https://edu.t-bio.info/course/transcriptomics-3/>
- <https://edu.t-bio.info/course/transcriptomics-4/>
- https://www.youtube.com/watch?v=WbJ9OA2vevk&feature=youtu.be&ab_channel=PineBiotech
- https://www.youtube.com/watch?v=UFB993xufUU&ab_channel=StatQuestwithJoshStarmer
- https://www.youtube.com/watch?v=Gi0JdrxRq5s&ab_channel=StatQuestwithJoshStarmer
- https://www.youtube.com/watch?v=tlf6wYJrwKY&ab_channel=StatQuestwithJoshStarmer



5. Practical

