

# TGAC MiRNA WORKSHOP

# DAY I MORNING SESSION I

DAY I MORNING SESSION I

# RAW DATA PROCESSING - QUALITY CONTROL - NORMALISATION - ANNOTATION

# SEQUENCE DATA FILE FORMATS

## FASTQ

- Raw sequencing information from the instrumentation and quality scores

```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
! ' ' * ( ( ( ***+ ) ) % % % ++ ) ( % % % % ) . 1 *** - + * ' ' ) ) ** 55CCF >>>> CCCCCCCC65
```

# SEQUENCE DATA FILE FORMATS

## FASTA

- Arbitrary header and processed sequences

```
>HEADER 1
TCCTCTCTCACAAAGCCAGGCTAT
>HEADER 2
TTTCAGATAATAGAAGTGAAGCGT
>HEADER 3
TCGGTTCCTGAATGGTTTT
```

# SEQUENCE DATA FILE FORMATS

## GFF

- Sequence Annotations (more on this later...)
- (up to) 9 Tab separated values

Chr1	TAIR10	CDS	30902	31079	.	+	1	Parent=AT1G01040.2,AT1G01040.2-Protein;
Chr1	TAIR10	three_prime_UTR	31080	31120	.	+	.	Parent=AT1G01040.2
Chr1	TAIR10	gene	28500	28706	.	+	.	ID=AT1G01046;Note=miRNA;Name=AT1G01046
Chr1	TAIR10	miRNA	28500	28706	.	+	.	ID=AT1G01046.1;Parent=AT1G01046;Name=AT1G01046.1;Index=1
Chr1	TAIR10	exon	28500	28706	.	+	.	Parent=AT1G01046.1
Chr1	TAIR10	gene	31170	33153	.	-	.	ID=AT1G01050;Note=protein_coding_gene;Name=AT1G01050
Chr1	TAIR10	mRNA	31170	33153	.	-	.	ID=AT1G01050.1;Parent=AT1G01050;Name=AT1G01050.1;Index=1
Chr1	TAIR10	protein	31382	32670	.	-	.	ID=AT1G01050.1-Protein;Name=AT1G01050.1;Derives_from=AT1G01050.1
Chr1	TAIR10	five_prime_UTR	33029	33153	.	-	.	Parent=AT1G01050.1
Chr1	TAIR10	exon	33029	33153	.	-	.	Parent=AT1G01050.1
Chr1	TAIR10	CDS	32547	32670	.	-	0	Parent=AT1G01050.1,AT1G01050.1-Protein;
Chr1	TAIR10	exon	32547	32670	.	-	.	Parent=AT1G01050.1
Chr1	TAIR10	CDS	32431	32459	.	-	2	Parent=AT1G01050.1,AT1G01050.1-Protein;
Chr1	TAIR10	exon	32431	32459	.	-	.	Parent=AT1G01050.1

# FASTQC

## INITIAL QUALITY CONTROL

- FASTQC
  - Most sequencers will give a basic overview of the run but this is not enough
- Check for technical quality of sequencing run
  - Looking for flow cell errors
  - machine errors
  - overall sequencing quality
- ‘Normal’ is of course dependant on what you might expect to see, certain samples will be bias for good reason

# FASTQC

- Basic Statistics
  - Basic information about total read counts
- Per base sequence quality
  - PHRED scores
- Per tile sequence quality
- Per sequence quality scores

# FASTQC

- Per base sequence content
- Per sequence GC content
- Per base N content
- Sequence Length Distribution

# FASTQC

- Sequence Duplication Levels
- Overrepresented sequences
- Adapter Content
- Kmer Content

# FASTQC EXERCISE

- Download the FASTQC program from the Babraham institute website
  - <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- Run the program by typing cd ~/Desktop/fastqc
  - Type perl fastqc
- Load the data found in ~/Desktop/data/FASTQ/SRR873382.fastq

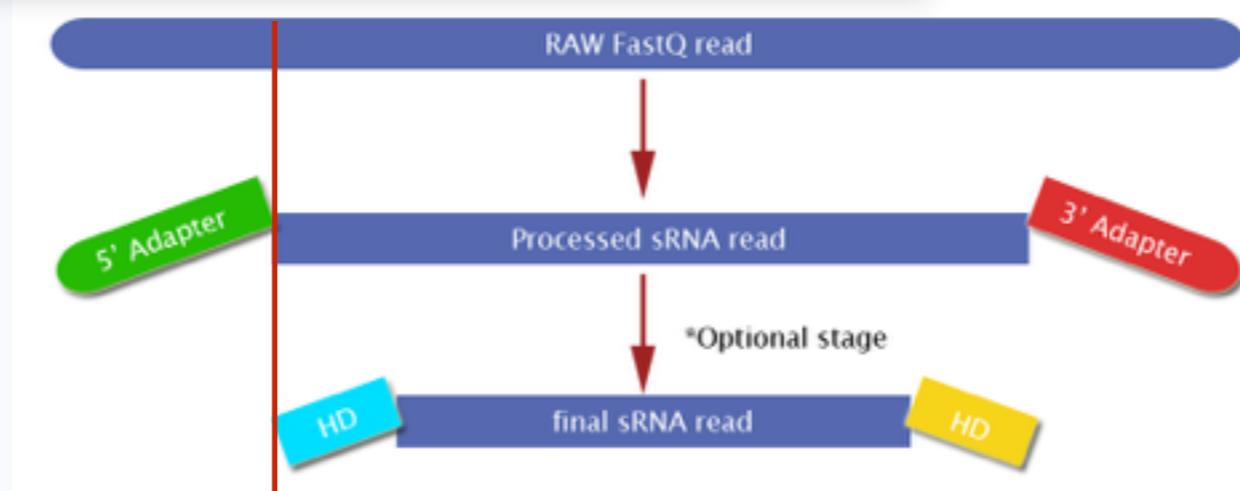
# WHAT CAN WE TELL?

- Per base sequence quality ok(ish)
- Per tile is not, why?
  - Is this a disaster?
- Per sequence quality scores are telling us that the majority of sequences are ok, but as we saw in the per tile, a problem occurred with a number of sequences
- Why is the nucleotide composition so spiky?

## WHAT CAN WE TELL?

- How about the Sequence Duplication level, FASTQC indicates a fail, but is this really a problem?
- KMER content, overrepresented motifs, what would we expect to see here?

# TRIMMING ADAPTERS



- sRNA sequences range from 18-40nt, however raw sequence reads from the illumina platform are generally selected as 50bp and therefore will potentially contain traces of the adapter sequence
- For the Illumina instruments these sequences will only be found at the 3' end of a read as the 5' end is used as a primer
- Adapter fragments must be removed prior to initial processing
- In addition to Illumina adapters certain protocols may introduce further inserts that require removal
- **HD adapter protocol** [Sorefan K et al. Reducing ligation bias of small RNAs in libraries for next generation sequencing. Silence. 2012;3: 4]

# SOFTWARE

## THE UEA sRNA WORKBENCH

- A suite of tools for various data preparation and analysis of sRNA-seq datasets
- Developed at the University of East Anglia and The Genome Analysis Centre
- Freely available from <http://srna-workbench.cmp.uea.ac.uk>
- Visit the web page and download the software
- Extract the file and place it on the Desktop

# SEQUENCE DATA PREPARATION

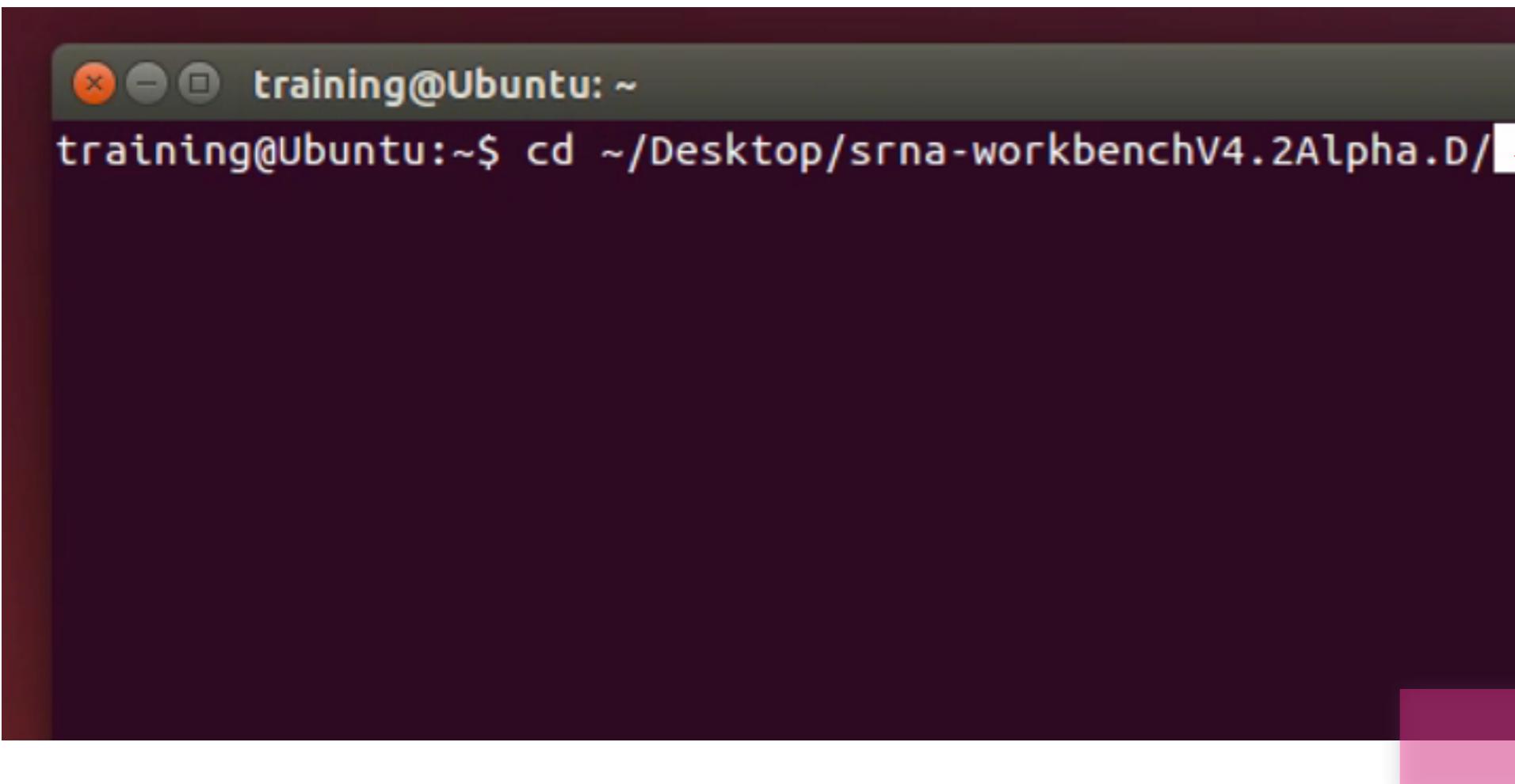
## TRIMMING ADAPTERS

- Using the sRNA Workbench exercise
- Adapter Sequence: ATCTCGTATGCCGTCT
- Sample investigating the effects of hypoxic conditions on MCF7 cells in H.Sapiens
- Load the Normoxia (time point 0) FASTQ sample

# ADAPTER REMOVAL EXERCISE

STRIP THE ADAPTERS AND FILTER

1. Load the Terminal App. Type: cd ~/Desktop/srna-workbenchV4.2.1Alpha.D



```
training@Ubuntu: ~
training@Ubuntu:~$ cd ~/Desktop/srna-workbenchV4.2.1Alpha.D/
```

# ADAPTER REMOVAL EXERCISE

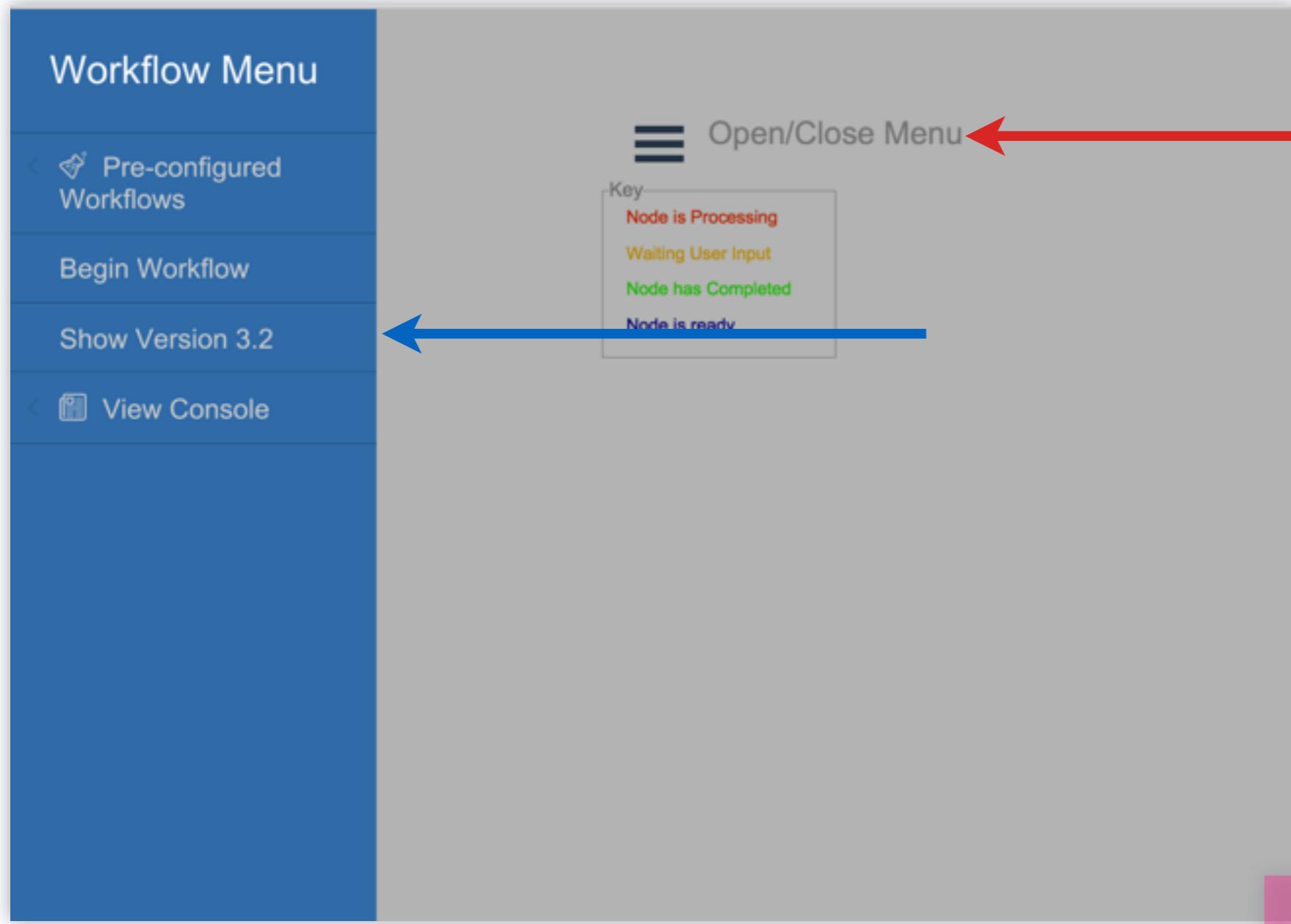
## RUN THE sRNA WORKBENCH

1. Type `java -Xms10g -Xmx10g -jar Workbench.jar`

```
training@Ubuntu: ~/Desktop/release
training@Ubuntu:~$ cd ~/Desktop/release/
training@Ubuntu:~/Desktop/release$ java -Xms10g -Xmx10g -jar Workbench.jar
```

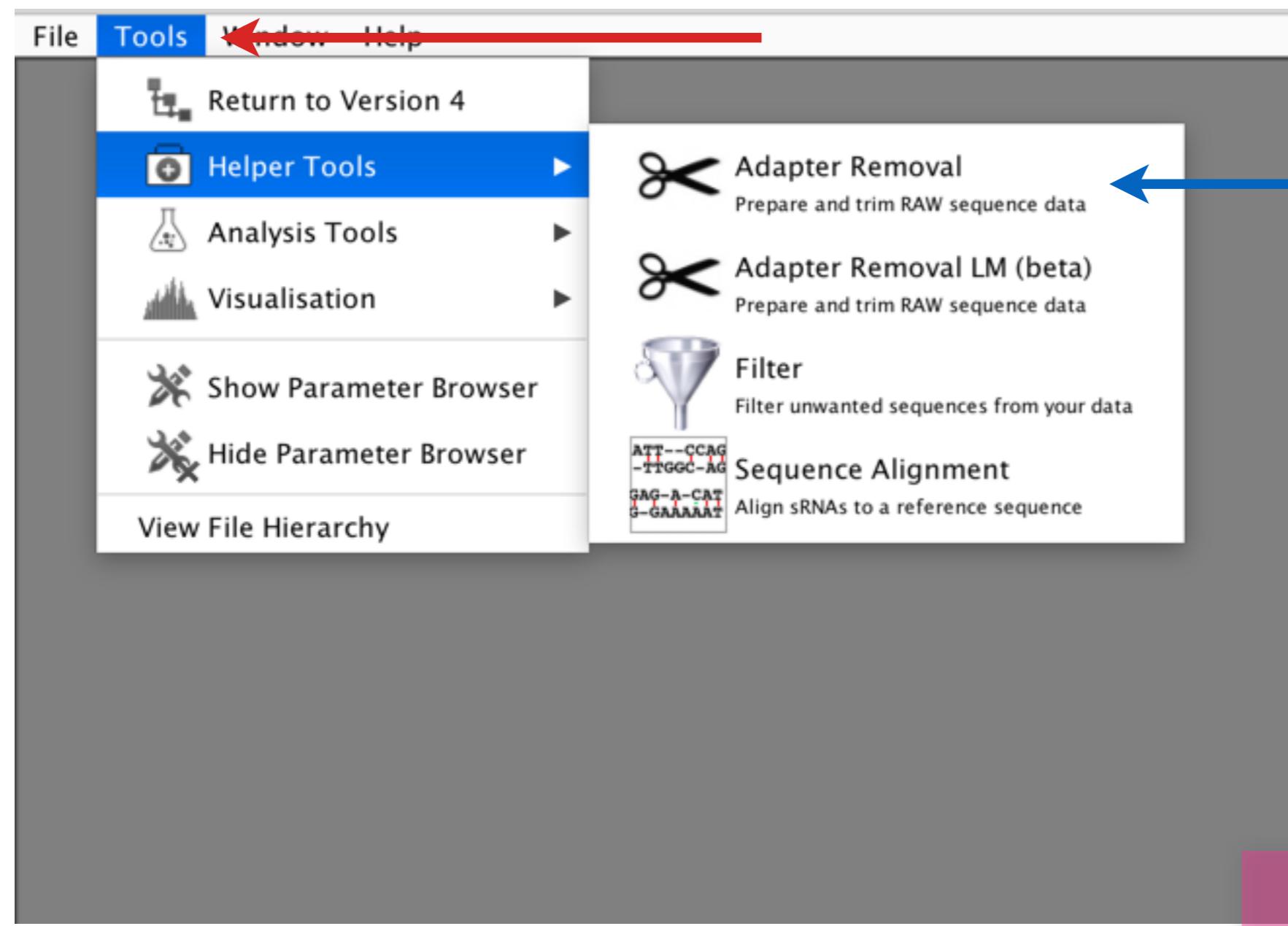
# ADAPTER REMOVAL EXERCISE

## LOAD THE V3.2 WORKBENCH



# ADAPTER REMOVAL EXERCISE

## LOAD THE ADAPTER REMOVAL TOOL



# REMOVING ADAPTERS



Adaptor Remover Options

File I/O

Input File Path (mandatory):  → ... D

Output File Directory (mandatory):  ... D

Discarded Sequence Directory (Optional):  ... D

Force Overwriting of files

Adaptor Selection

HD Adaptor Protocol

Enable HD Adaptor Processing HD Full HD Simple 5' Prime HD Simple 3' Prime

3' Adaptor (mandatory)

Adaptor Sequence:

Pre-defined Adaptors:  ← ...

Nucleotides to Use:  (min. 5, max. 20)

5' Adaptor (optional)

Adaptor Sequence:

Pre-defined Adaptors:  ← ...

Nucleotides to Use:  (min. 5, max. 20)

Filtering

Min. Length:  (min. 1) Max. Length:  (max. 80)

Open:  
~/Desktop/data/FASTQ  
/SRR873382.fastq

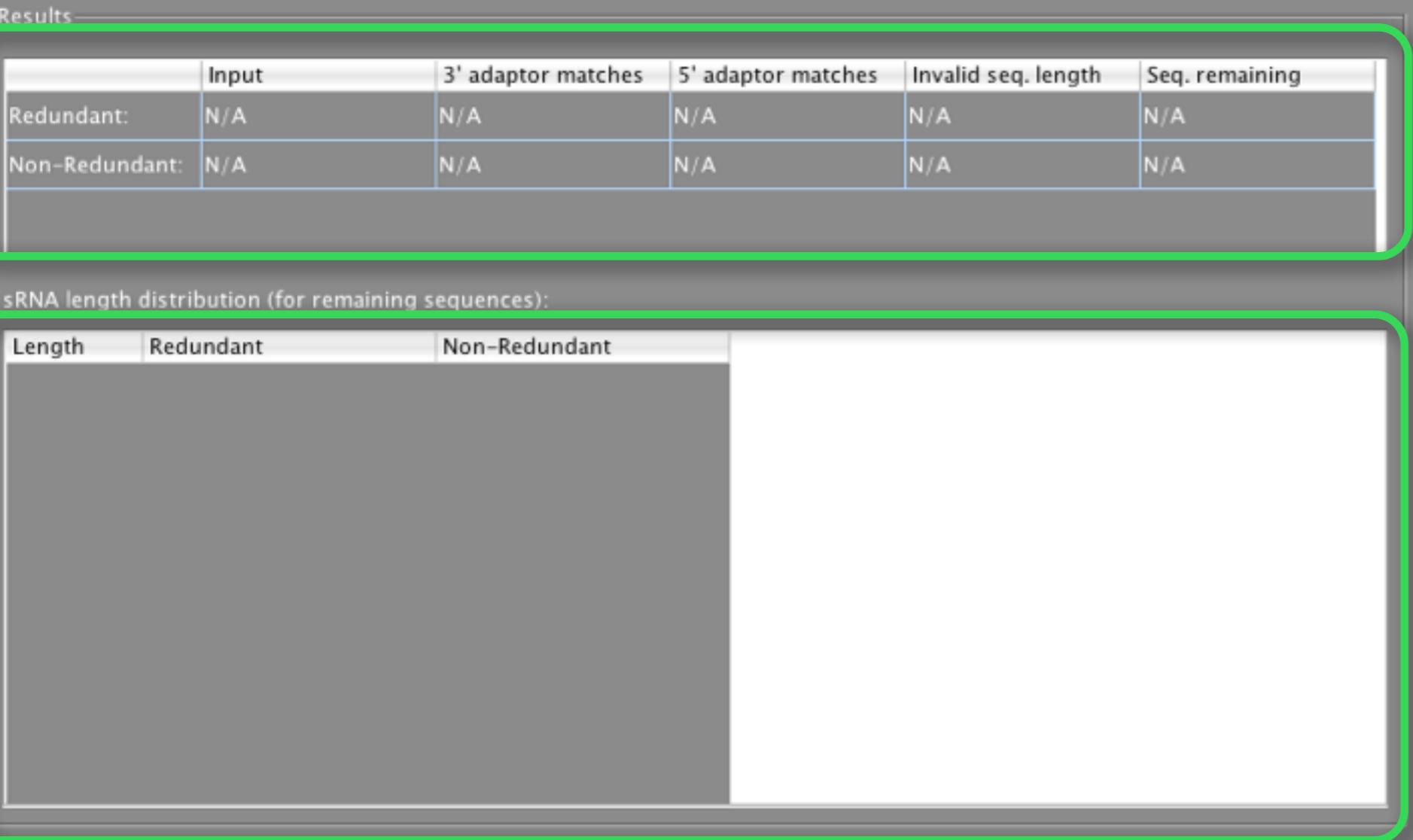
Output:  
~/Desktop/data/results

Select the LMN2 3'  
(Illumina Tru Seq 1.5)  
Adapter sequence

# RESULTS



## Process Overview



## Length Distributions

# FURTHER QUALITY CHECKING

## SETUP PIPELINE

### LOAD THE QUALITY CHECK PIPELINE

1. Load the Terminal App. Type: 'cd ~/Desktop/srna-workbench
2. Type java -Xms10g -Xmx10g -jar Workbench.jar

# SETUP PIPELINE

## LOAD THE QC/DE PIPELINE

### Workflow Menu

<  Pre-configured Workflows

Begin Workflow

Show Version 3.2

<  View Console

BACK

Create Quality Check Workflow

Create Quality Check and Differential Expression Workflow

# SETUP PIPELINE

LOAD THE WIZARD



# SETUP PIPELINE

## Sequence Data Setup

Show Wizard

Database Settings

Export

Load previous

Home

## CONFIGURE THE SAMPLES

### 1 Step 1 Setup Samples

### 2 Step 2 Select References

### 3 Step 3 Configure Annotations

#### Sample Setup

What is the ID for this sample?

Add Sample

Sample ID	Small RNA Replicate Details	
1 ×	Replicate Number	Filename
<a href="#">Add Files</a>		

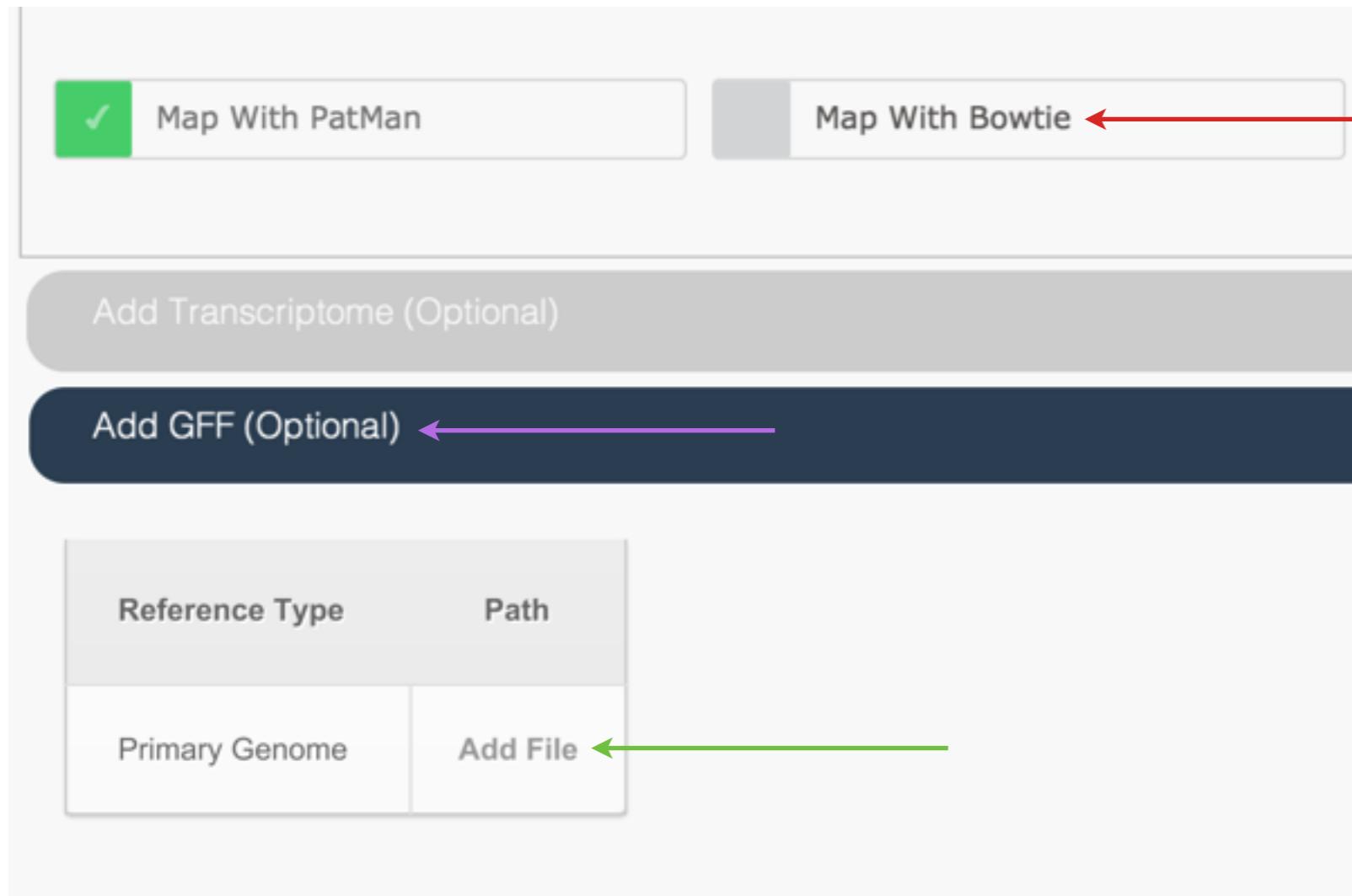
Load the files found in ~Desktop/data/FASTA

Sample 1: SRR873382\_na\_sampled\_0.1.fasta, SRR873383\_na\_sampled\_0.1.fasta

Sample 2: SRR873384\_na\_sampled\_0.1.fasta, SRR873385\_na\_sampled\_0.1.fasta

# SETUP PIPELINE

## CONFIGURE THE REFERENCES



Select “Map with Bowtie”

Load the genome file found in ~Desktop/data/FASTA/genome/GRCh38\_no\_alt

and the GFF file found in ~/Desktop/data/GFF

# SETUP PIPELINE

## 1 Step 1

Setup Samples

## 2 Step 2

Select References

## 3 Step 3

Configure Annotations

Available Annotation	Selected Annotation	Add To 'Other'
	miRNA_primary_transcript miRNA	

Drag the miRNA and annotation into the selected options

# BEGIN PIPELINE

## Sequence Data Setup

Show Wizard

< Database Settings

< Export

Load previous

Home

## Workflow Menu

< Pre-configured Workflows

Begin Workflow

Show Version 3.2

< View Console



# QUALITY CONTROL

- An often overlooked step...
- Checking the sequencing information for issues that may impact downstream analysis
  - Fails: Technical
    - Library preparation
  - Fails: Biological
    - DNA or other contamination to RNA?
    - Sample infection?
- Various indications on the quality of both biological sample and sequencing run can be made through a set of diagnostic plots
- Summarise, Visualise but not necessarily discarding!



SWITCH BACK TO SOFTWARE

## QUALITY CONTROL

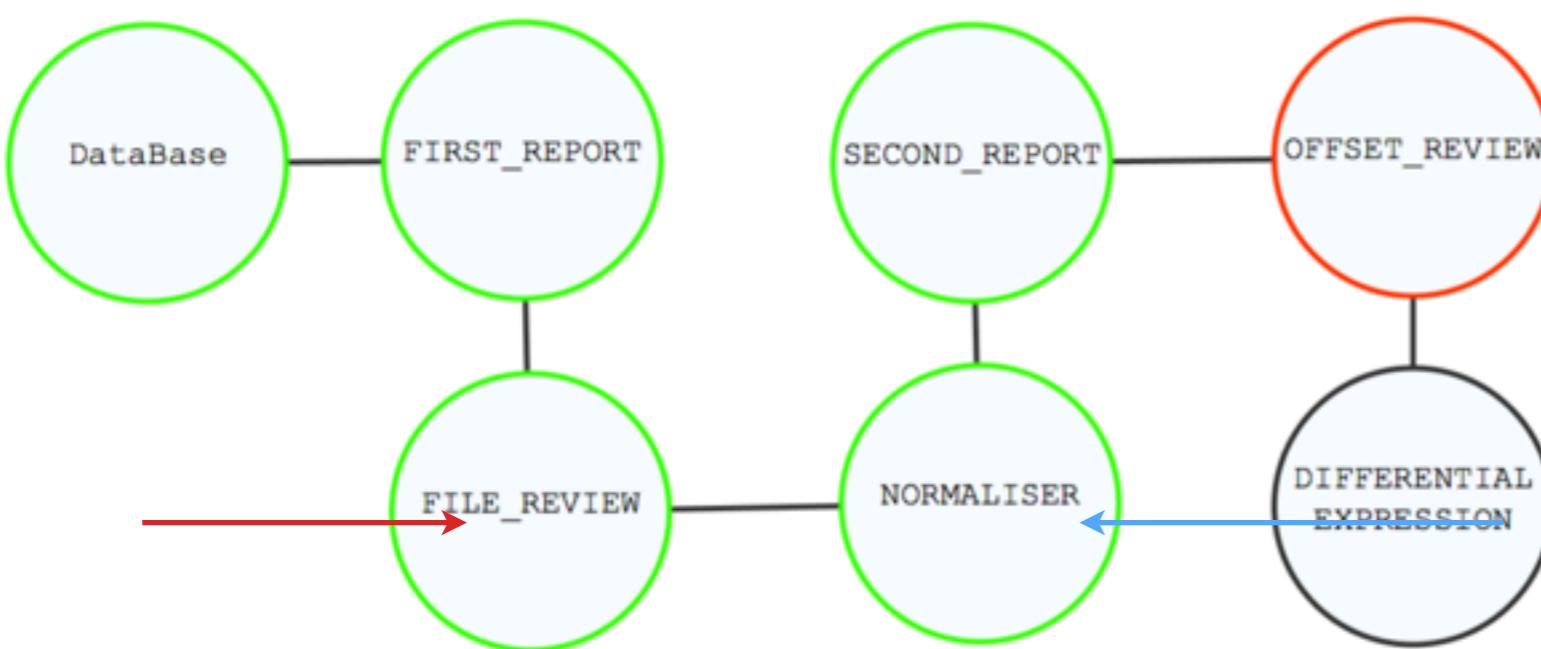
- Investigate the plots found in the first report node of the workflow
- Do the plots indicate comparable sets of replicates and samples?
- How about the mapping quality?

# NORMALISATION

# CONTINUE PIPELINE

## REVIEW REPLICATION QUALITY

1. Go into the file review node, and hit continue
2. Continue to Normalisation
3. Select All Normalisations and Continue



# CONTINUE PIPELINE

## File Review

Continue Workflow

Home

Size Class Controls



Open/Close Menu

Size Class Filtering

17

19



18

20

File Review and Removal

Sample ID

Small RNA Replicate Details

Degradome Replicate De

Hierarchy Visualisation

# CONTINUE PIPELINE

Normalisation Setup

Export

Continue Workflow

Home

Open/Close Menu

Scaling Normalisation

Per Total Normalisation    Upper Quartile Normalisation    Trimmed Mean Normalisation    DESEQ Normalisation

Rank Based Normalisation

Quantile Normalisation

Statistical Normalisation

Bootstrapping Normalisation



# NORMALISATION

## WHY DO WE DO IT?

- To attempt to correct issues found during the initial quality control investigation such as highly differentially expressed reads swallowing up the sequencing space
- Sequencing instruments do not have pre-determined sequencing depths
- When sequencing multiple samples expression levels should be adjusted to take this into account
- So if attempting to identify differentially expressed reads over a set of reads the false positive rate is reduced

# NORMALISATION

## THREE DISTINCT CATEGORIES

- Scaling
  - RPM (Reads per million\*)
  - TMM Trimmed mean of means (edgeR)
  - DESeq 2
- Rank based
  - Quantile
- Statistical
  - Subsampling and Bootstrapping

## NORMALISATION QC

### WHY LOOK AT QUALITY CONTROL AGAIN?

- No single normalisation is suitable for all datasets
- Return to the issues that may have been discovered during QC and investigate the effect made by each normalisation



## SWITCH BACK TO SOFTWARE

## NORMALISATION QC

- Investigate the box plots, fold change box plots and MA plot by normalisation, which normalisation is suitable for this data?
-

# SEQUENCE ANNOTATION

# ANNOTATION

## WHAT DID WE SEQUENCE?

- Several Sources of annotations
  - GFF
  - BLAST
  - RFAM
  - miRBase

# ANNOTATED QUALITY REPORTS

# GFF ANNOTATION

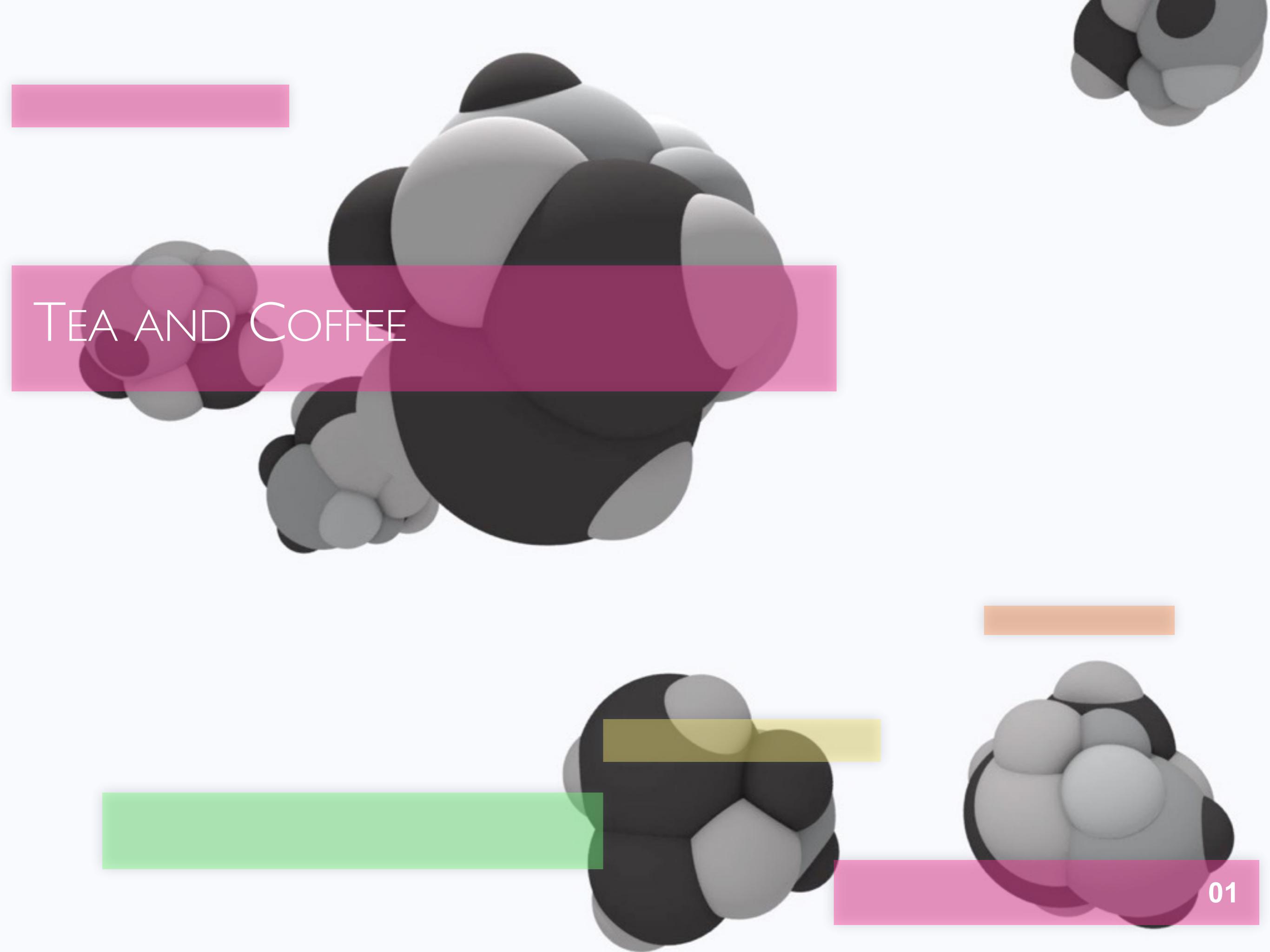
- Align reads to a ‘known’ set of sRNAs
- Annotation information can be used in later analysis for example in differential expression analysis
- Certain Quality Checks should exhibit certain properties
  - miRNA should show a tight cone distribution between replicates
  - t/r RNA are expected to show dispersed distribution between replicates (large amounts of degradation)
- Requires all possible\* alignments of an sRNA sequence



## RETURNING TO THE REPORT

## ANNOTATED QC

- From the plots, is the proportion of non-redundant miRNAs that map to the genome what you may expect from this dataset?
  - If something has gone very wrong with the data, it may be worth looking at this proportion for other classes of sRNA such as t and rRNA
- Have a look at the MA plots for miRNAs, what shape would you expect?



# TEA AND COFFEE

# DAY I MORNING SESSION 2

# DAY I MORNING SESSION 2

## MiRNA ANALYSIS

- Overview of miRNAs
- Annotating known miRNAs from miRBase
- Profiling miRNAs in deep sequencing data
- How to analyse miRNA prediction data
- How to spot false positives

# MiRBASE ANNOTATION

# MIRBASE ANNOTATION

- <http://www.mirbase.org/>
- Central repository of annotated miRNAs
- Contains all (or the vast majority) of published miRNAs
- Provides a search function
- Provides direct experimental evidence for miRNA expression from NGS data
- Is downloadable and machine parsable

# MIRBASE ANNOTATION

Stem-loop sequence hsa-mir-23a	
<b>Accession</b>	MI0000079
<b>Symbol</b>	<a href="#">HGNC:MIR23A</a>
<b>Description</b>	Homo sapiens miR-23a stem-loop
<b>Gene family</b>	MIPF0000027; <a href="#">mir-23</a>
<b>Stem-loop</b>	 <a href="#">Get sequence</a>
<b>Deep sequencing</b>	<p><a href="#">89501</a> reads, 57 experiments</p>  <p># reads</p> <p>GGCGGGCUGGGGUUCCUGGGGAUGGGAUUJGCUUCCUGUCACAAUCAUUGCCAGGGAUUCCAACCGACC</p>
<b>Comments</b>	This miRNA was previously named miR-23 [1,2] but is renamed here to avoid confusion with the more recently described miR-23b ( <a href="#">MI0000439</a> ). Kawasaki and Taira reported that miR-23 regulates the transcriptional repressor Hairy enhancer of split (HES1) [3]. This finding was later retracted after the discovery that the regulated gene was human homolog of ES1 (HES1), whose function is unknown.
<b>Genome context</b>	<p><i>Coordinates (GRCh37)</i> <a href="#">19: 13947401-13947473 [-]</a></p> <p><i>Overlapping transcripts</i> intergenic</p>
<b>Clustered miRNAs</b>	<p>&lt; 10kb from hsa-mir-23a</p> <p>hsa-mir-23a      <a href="#">19: 13947401-13947473 [-]</a></p> <p><a href="#">hsa-mir-27a</a>      <a href="#">19: 13947254-13947331 [-]</a></p> <p><a href="#">hsa-mir-24-2</a>      <a href="#">19: 13947101-13947173 [-]</a></p>

# MIRBASE ANNOTATION

## Mature sequence hsa-miR-23a-5p

Accession	<a href="#">MIMAT0004496</a>
Previous IDs	hsa-miR-23a*
Sequence	9 - <chem>gggguucuggggauuu</chem> - 30 <a href="#">Get sequence</a>
Deep sequencing	<a href="#">330</a> reads, 19 experiments
Evidence	experimental; cloned [6-7]
Predicted targets	DIANA-MICROT: <a href="#">hsa-miR-23a-5p</a> MICRORNA.ORG: <a href="#">hsa-miR-23a-5p</a> MIRDB: <a href="#">hsa-miR-23a-5p</a> PICTAR-VERT: <a href="#">hsa-miR-23a</a>

## Mature sequence hsa-miR-23a-3p

Accession	<a href="#">MIMAT0000078</a>
Previous IDs	hsa-miR-23a
Sequence	45 - <chem>aucacauugccaggauucc</chem> - 65 <a href="#">Get sequence</a>
Deep sequencing	<a href="#">89229</a> reads, 52 experiments
Evidence	experimental; cloned [1,4-7], Northern [1]
Validated targets	TARBASE: <a href="#">hsa-miR-23a-3p</a>
Predicted targets	DIANA-MICROT: <a href="#">hsa-miR-23a-3p</a> MICRORNA.ORG: <a href="#">hsa-miR-23a-3p</a> MIRDB: <a href="#">hsa-miR-23a-3p</a> TARGETSCAN-VERT: <a href="#">hsa-miR-23a</a> PICTAR-VERT: <a href="#">hsa-miR-23a</a>

# MIRBASE ANNOTATION

## References

- 1 PMID:[11679670](#)  
["Identification of novel genes coding for small expressed RNAs"](#)  
Lagos-Quintana M, Rauhut R, Lendeckel W, Tuschl T  
Science. 294:853-858(2001).
- 2 PMID:[11914277](#)  
["miRNPs: a novel class of ribonucleoproteins containing numerous microRNAs"](#)  
Mourelatos Z, Dostie J, Paushkin S, Sharma A, Charroux B, Abel L, Rappaport J, Mann M, Dreyfuss G  
Genes Dev. 16:720-728(2002).
- 3 PMID:[12808467](#)  
["Hes1 is a target of microRNA-23 during retinoic-acid-induced neuronal differentiation of NT2 cells"](#)  
Kawasaki H, Taira K  
Nature. 423:838-842(2003).
- 4 PMID:[14573789](#)  
["Reduced accumulation of specific microRNAs in colorectal neoplasia"](#)  
Michael MZ, O' Connor SM, van Holst Pellekaan NG, Young GP, James RJ  
Mol Cancer Res. 1:882-891(2003).
- 5 PMID:[15325244](#)  
["Altered expression profiles of microRNAs during TPA-induced differentiation of HL-60 cells"](#)  
Kasashima K, Nakamura Y, Kozu T  
Biochem Biophys Res Commun. 322:403-410(2004).
- 6 PMID:[17604727](#)  
["A mammalian microRNA expression atlas based on small RNA library sequencing"](#)  
Landgraf P, Rusu M, Sheridan R, Sewer A, Iovino N, Aravin A, Pfeffer S, Rice A, Kamphorst AO, Landthaler M, Lin C, Socci ND, Hermida L, Fulci V, Chiaretti S, Foa R, Schliwka J, Fuchs U, Novosel A, Muller RU, Schermer B, Bissels U, Inman J, Phan Q, Chien M  
Cell. 129:1401-1414(2007).
- 7 PMID:[17616659](#)  
["Patterns of known and novel small RNAs in human cervical cancer"](#)  
Lui WO, Pourmand N, Patterson BK, Fire A  
Cancer Res. 67:6031-6043(2007).

Deep sequencing reads for stem-loop sequence MI0000079

Stem-loop ID	Reads	hsa-mir-23a	hsa-mir-23a-5p	hsa-mir-23a-3p	Count	RPM
			<u>GGGGGUUCCUGGGGAUGGGAUU</u>		29	0.332
			<u>GGGGGUUCCUGGGGAUGGGAU</u>		10	0.0923
			<u>GGGGGUUCCUGGGGAUGGGAUU</u>		215	2.81
			<u>GGGGGUUCCUGGGGAUGGGAUU</u>		34	0.316
			<u>GGGGGUUCCUGGGGAUGGGAU</u>		10	1.97
				<u>AAAUCACAUGCCAGGGAUUUCA</u>	403	678
				<u>AAAUCACAUGCCAGGGAUUUCC</u>	348	708
				<u>AAAUCACAUGCCAGGGAUUUCAA</u>	77	103
				<u>AAAUCACAUGCCAGGGAUUUCAAAC</u>	75	106
				<u>AAAUCACAUGCCAGGGAUUUUC</u>	67	166
				<u>AAAUCACAUGCCAGGGAUUU</u>	46	107
				<u>AAAUCACAUGCCAGGGAUU</u>	23	79.1
				<u>AAUCACAUGCCAGGGAUUUCA</u>	148	2.33
				<u>AAUCACAUGCCAGGGAUUUCC</u>	86	60.5
				<u>AAUCACAUGCCAGGGAUUU</u>	69	27.6
				<u>AAUCACAUGCCAGGGAUUUC</u>	45	11.8
				<u>AAUCACAUGCCAGGGAUUUCAA</u>	23	0.265
				<u>AUCACAUGCCAGGGAUUUCA</u>	28827	404
				<u>AUCACAUGCCAGGGAUUUCAA</u>	13214	164
				<u>AUCACAUGCCAGGGAUUUCC</u>	6167	2.29e+03
				<u>AUCACAUGCCAGGGAUUUC</u>	5293	541
				<u>AUCACAUGCCAGGGAUUU</u>	4248	783
				<u>AUCACAUGCCAGGGAUU</u>	838	249
				<u>AUCACAUGCCAGGGAU</u>	129	45.1
				<u>AUCACAUGCCAGGGAUUUCAAAC</u>	46	0.479
				<u>AUCACAUGCCAGGG</u>	35	46.3
				<u>UCACAUGCCAGGGAUUUCAA</u>	128	2.04
				<u>UCACAUGCCAGGGAUUUCCA</u>	118	4.77
				<u>UCACAUGCCAGGGAUUUCAAAC</u>	62	1.76
				<u>UCACAUGCCAGGGAUUUCC</u>	48	22.1
				<u>ACAUGCCAGGGAUUUCCA</u>	66	0.748
				<u>ACAUGCCAGGGAUUUCAA</u>	48	0.568
				<u>ACAUGCCAGGGAUUUCC</u>	21	0.26
				<u>ACAUGCCAGGGAUUU</u>	13	0.217
				<u>ACAUGCCAGGGAUUUC</u>	13	0.211
		GGCCGGCU	<b>GGGGGUUCCUGGGGAUGGGAUUU</b>	<b>GCUUCCUGUCACAA</b>	<b>AUCACAUGCCAGGGAUUUCC</b>	<b>AACCGACC</b>
		(..(.	(.(.	(.(.	(.(.	(.-33.20)

# MIRBASE EXERCISE I

- Go to miRBase <http://www.mirbase.org/> and have a look at the kinds of data that miRBase contains.
- Find the miRBase entry for mir-100 in mouse.
- What other microRNAs are located close to mir-100 in mouse?
- What other microRNAs are related to mir-100?
- In what range of species are mir-100 family members found?

# MIRBASE EXERCISE 2

- Find mouse mir-196-2
  - Look at the mature sequences that derive from the mir-196-2 hairpin.
  - Have a look at the predicted targets of miR-196-2 sequences
  - Are there any validated target sites?
  - Click the genomic coordinates of the miRNA to link to the EnsEMBL genome annotation database.
  - Is there any genomic relationship between the miRNA and its validated targets?

Deep sequencing reads for stem-loop sequence MI0000079



# MIRBASE EXERCISE 3

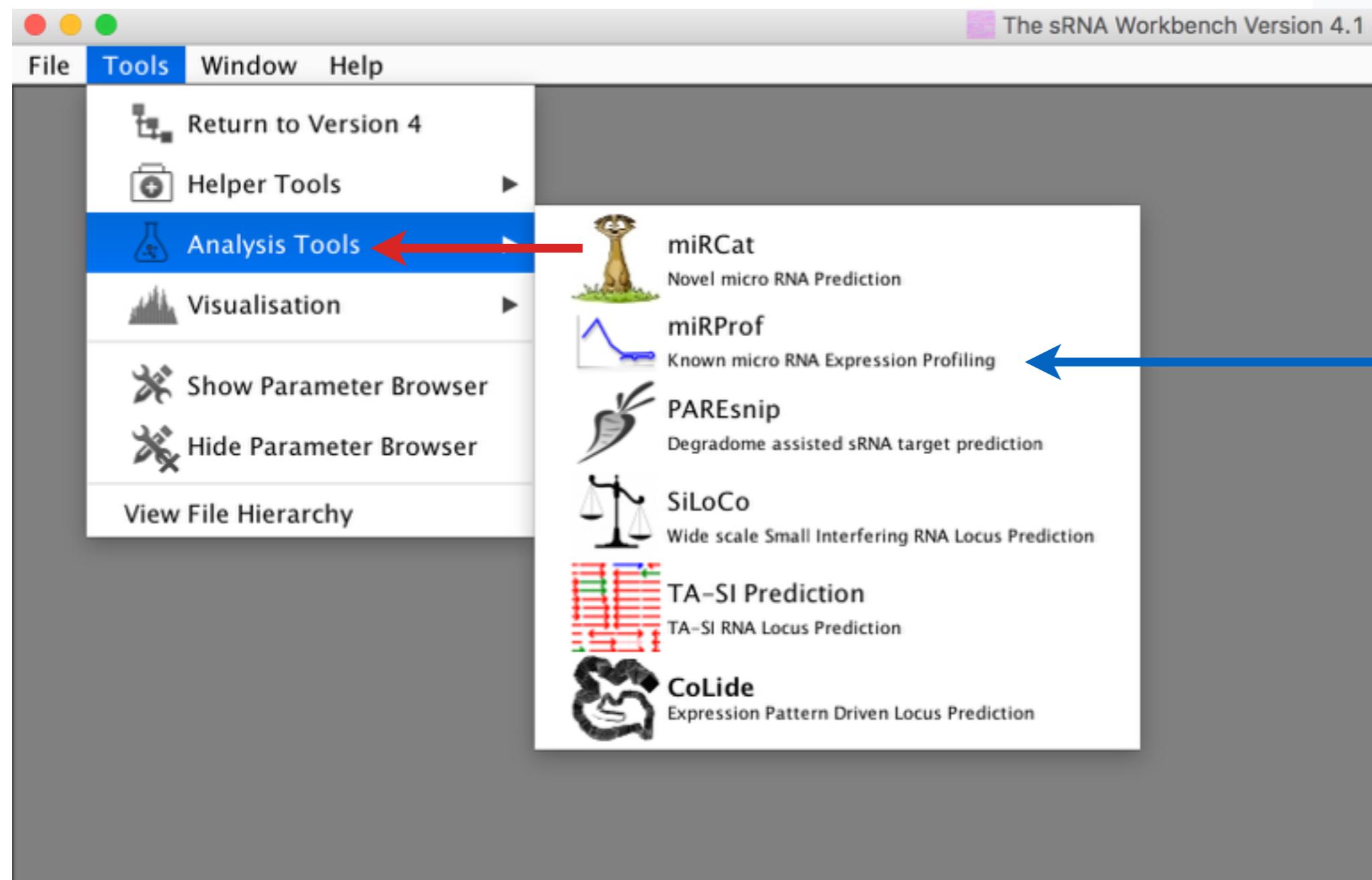
- Find the miRBase entry for miR-317 in *Drosophila melanogaster*.
  - Look at the pattern of reads that map to the locus from deep-sequencing experiments. Do the read counts support the mature sequence annotations?
  - Which mature sequence is most abundant?
- Now search for the microRNA osa-MIR439c.
  - Explore the expression pattern of the mature miRNAs given by the NGS data (how many reads? what is the most abundant read? etc.)
  - Considering the miRNA annotation guidelines, what can you say about this miRBase entry?
  - Look at other miRNAs in the same gene family - are they likely to be real miRNAs?
  - How many other microRNA were discovered in the same study?

# MIRBASE EXERCISE 4

- Use the browse and search pages for this exercise.
  - From the browse page, how many microRNAs are annotated in *Drosophila melanogaster*?
  - From the "tissue expression" section of the search page, find the list of all deep sequencing experiments in *D. melanogaster*.
  - Compare the miRNAs found in a head dataset with an embryonic dataset.
  - Identify miRNAs that are highly expressed in the head, but not in the embryo.

# MIRPROF EXERCISE

## LOAD THE MIRPROF TOOL

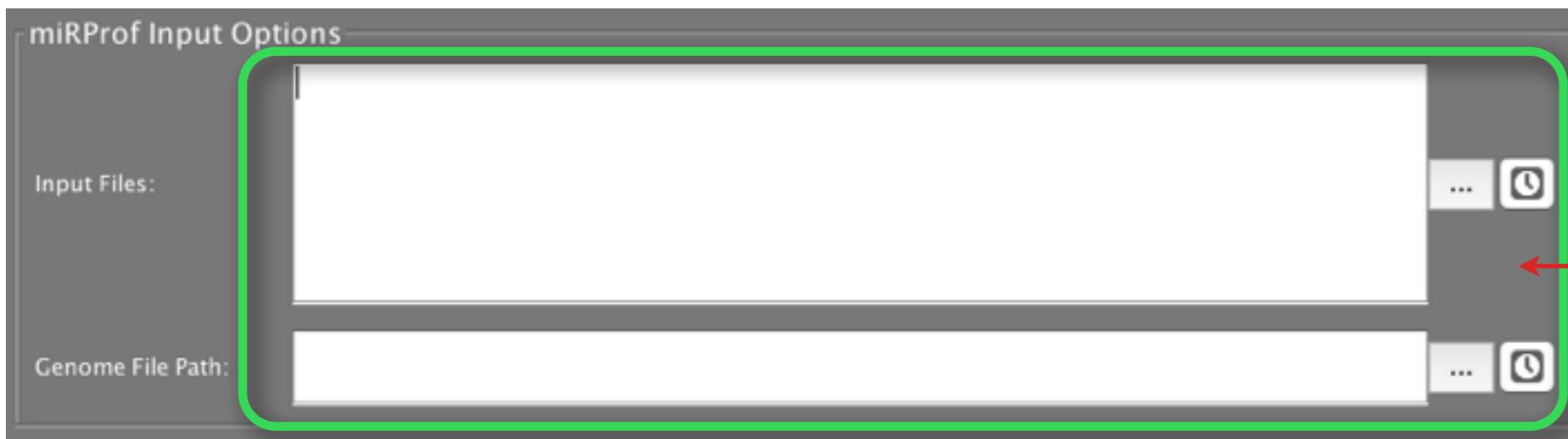


# MIRPROF EXERCISE I



- Visit the UCSC Genome Browser website
  - <http://genome.ucsc.edu/>
- Locate Chromosome 1 of the human genome and download it, then extract it to any location

# MIRPROF EXERCISE I



Add the files found in ~/Desktop/data/FASTA/

1.SRR873382\_na\_sampled\_0.1.fasta

2.SRR873384\_na\_sampled\_0.1.fasta

and the single chromosome file you downloaded

# MIRPROF EXERCISE



- Use the default parameters and run the program
- Export the results to csv and open them in the spreadsheet editor
- Calculate the simple fold change for each miRNA using the total counts
  - $(\text{Expression Sample 2} - \text{Expression Sample 1}) / \text{Expression Sample 1}$
- Are there differences in the miRNA Expression for each sample

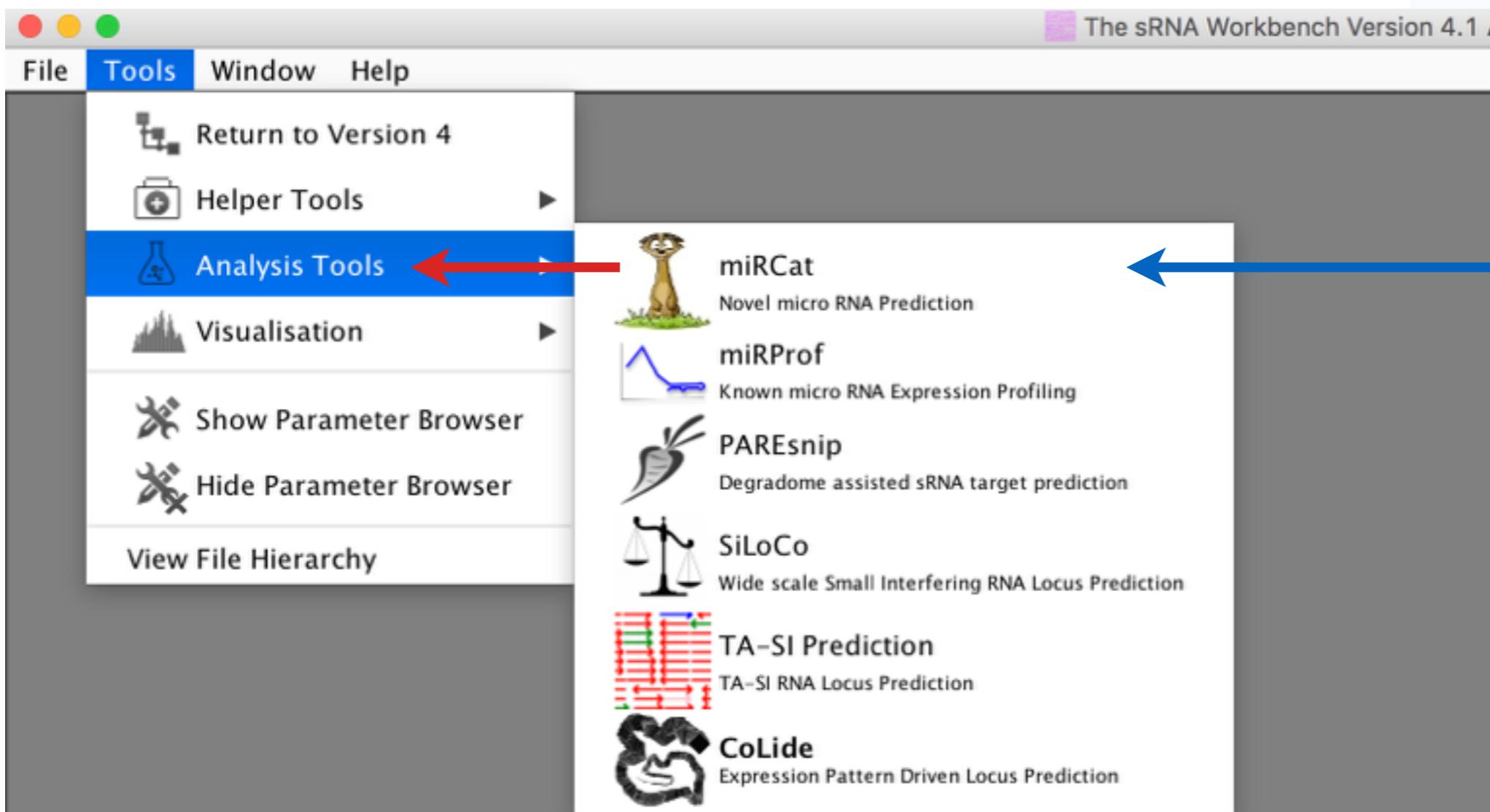
# MIRCAT EXERCISE



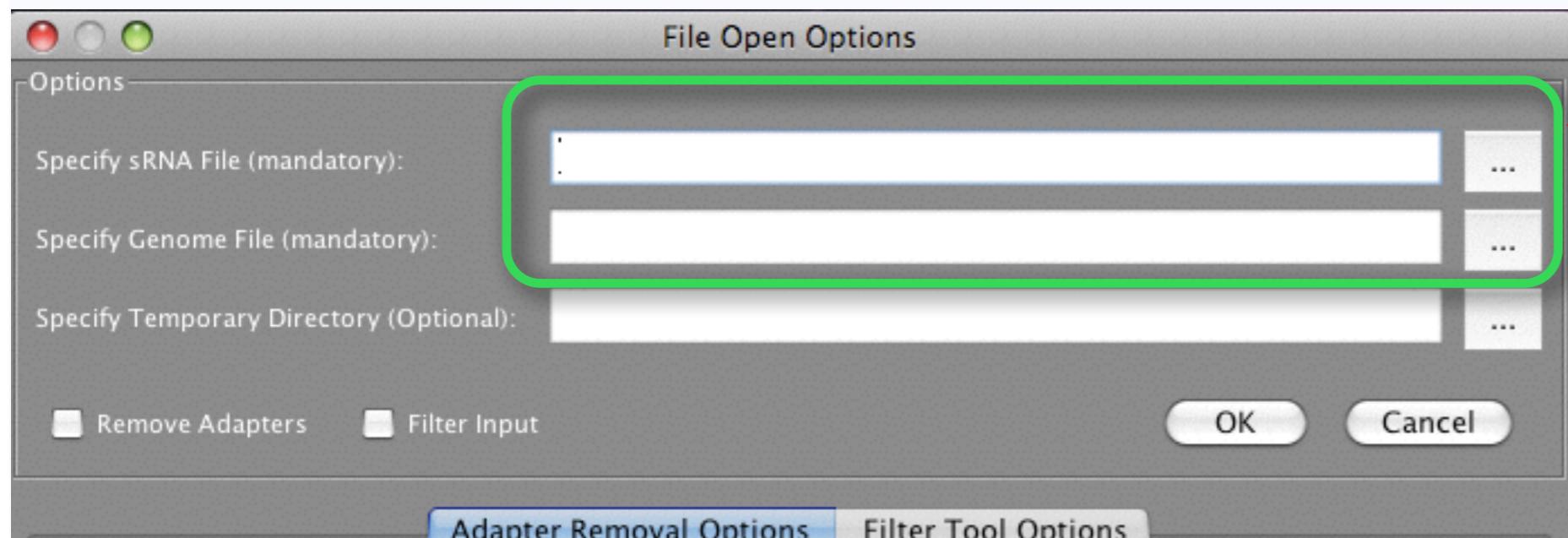
- Trim the adapters from the reads.fa file found in ~/Desktop/miRDeep2/tutorial\_dir

# MiRCAT EXERCISE

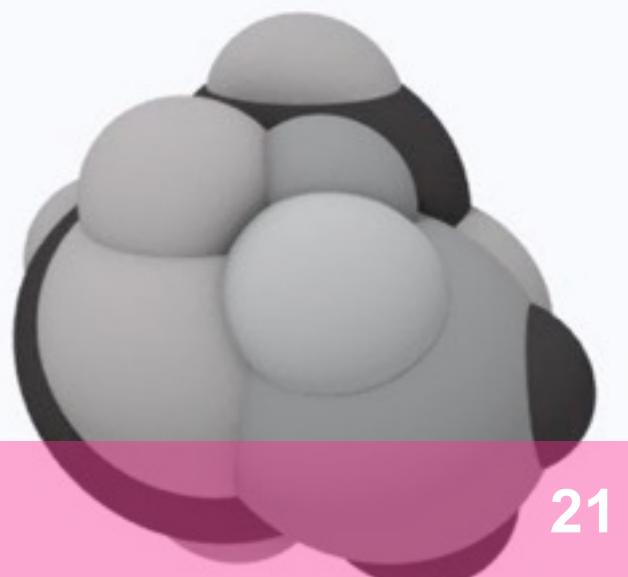
## LOAD THE MiRCAT TOOL



# MIRCAT EXERCISE



- Start a new project from the menu
- load the trimmed file you created into the sRNA input box
- and the genome file found in the original location called cel\_cluster.fa
- Use the default animal parameters and run the program



# MIRCAT EXERCISE



## Options

miRcat

File Run Help

Export Results to CSV Export miRNAs to FASTA Output Hairpins Render all Hairpins Cancel current run Show All Results In SeQViSsr

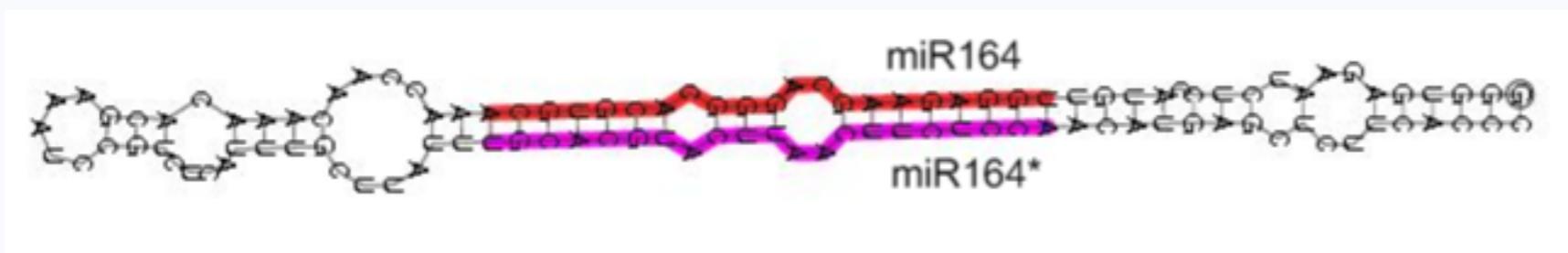
Chromosome	Sequence	miRNA*	Hairpin Sequence	miRNA Start Coordinate	miRNA End Coordinate
1	TGCAAGAAGGGAGAAGCAAAGT	TTTCTTCTACTTCTTGCACA(1)	GAGGACTTACATGGCCTCAAGTCACCTGTGGTGTG TGCAAGAAGGGAGAAGCAAAGT CTC TCTATGTATTATGAGATAGCTACTTCTATGGCTAGGATATATGTTGACAAGACCGGC TCTTCTACTTCTTGCACA ACCTGAGGTTATTGAGGCTATAAGTCTTC	28553	28573
1	TCGGACCAGGCTTCATCCCC	GGAATGTTGCTGGATCGAGGATA(1) GGAATGTTGCTGGATCGAGGA(3) GGAATGTTGCTGGATCGAGG(78)	ATGAGGGTTAACGCTATTCAAGTGAGG GGAATGTTGCTGGATCGAGG ATATTATAGAT ATATACATGTGATGTTAATGATTCAAGTGATCATAGAGAGTATCC TCGGACCAGGCTTC ATCCCCCCCCAACATGTTATTGCCCTGTGATCACCAT	78932	78952

Output table

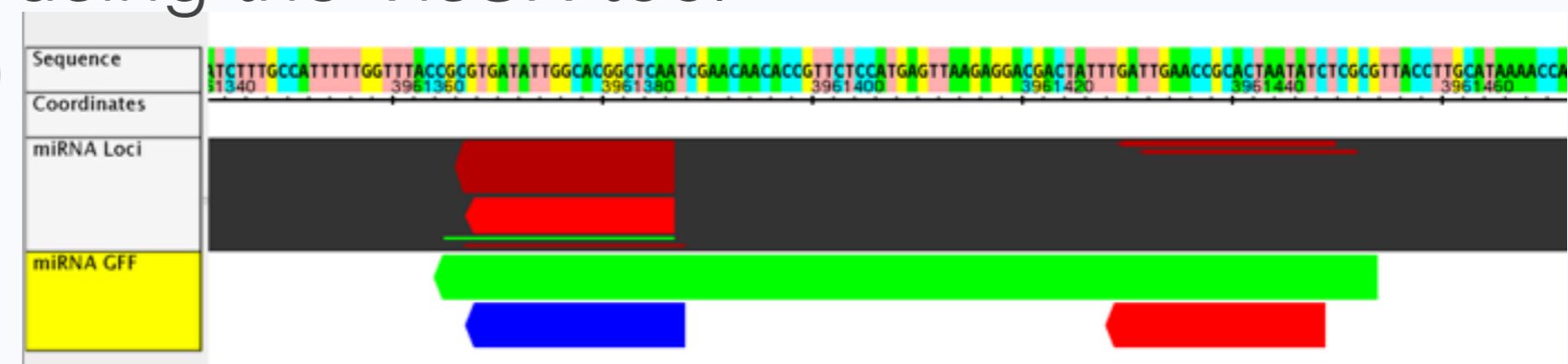
# MIRCAT EXERCISE I



- Pick one and generate the pre cursor plot for it (right click menu)



- Look at the locus using the VisSR tool (right click menu)



- Export all of the results to CSV file, we will use them later

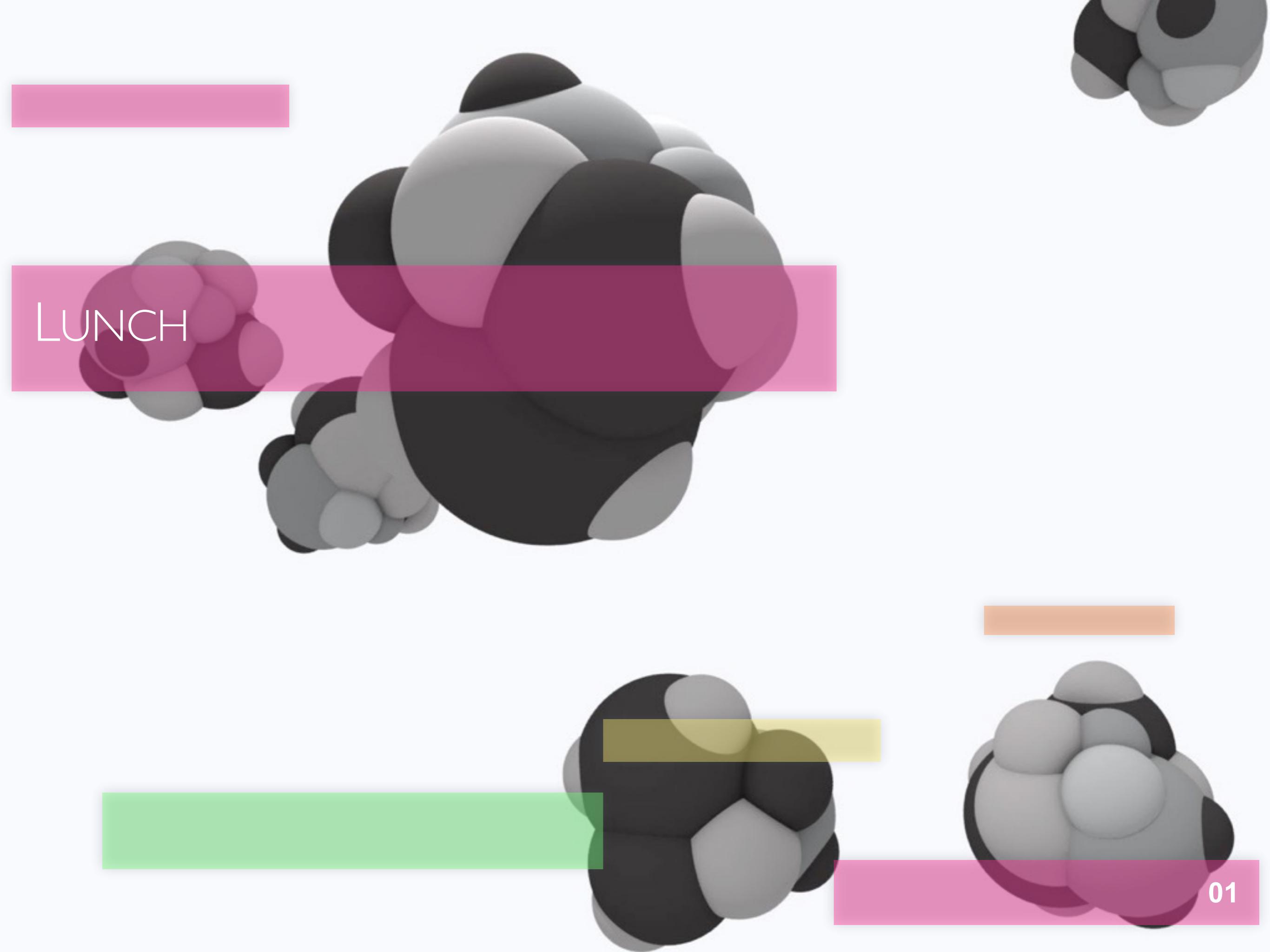
# MIRDEEP 2 EXERCISE I



# MiRCAT EXERCISE 2



- Examine the results from both the miRCat and miRDeep runs
- Are there any differences between the two?
- What can we infer from this?



LUNCH

# DAY I AFTERNOON SESSION I

# DIFFERENTIAL EXPRESSION

## DE EXERCISE I

- From the second report, with a fill dataset we would normally select the normalisation that best suits the data
- Double click in the Offset Review Node and choose the bootstrapping normalisation from the menu
- Select continue workflow

## DE EXERCISE I

- Look at the KL graph produced from the dataset
  - Does the curve line up with the global minimum for average abundances across the dataset?

# DE EXERCISE I

- Open the Differential Expression node.  
Compare Sample 1 with Sample 2 by  
dragging their references into the reference  
and observed columns of the interface
- Use the options menu to set the fold  
change cut off to 0.5
- Investigate the result set