

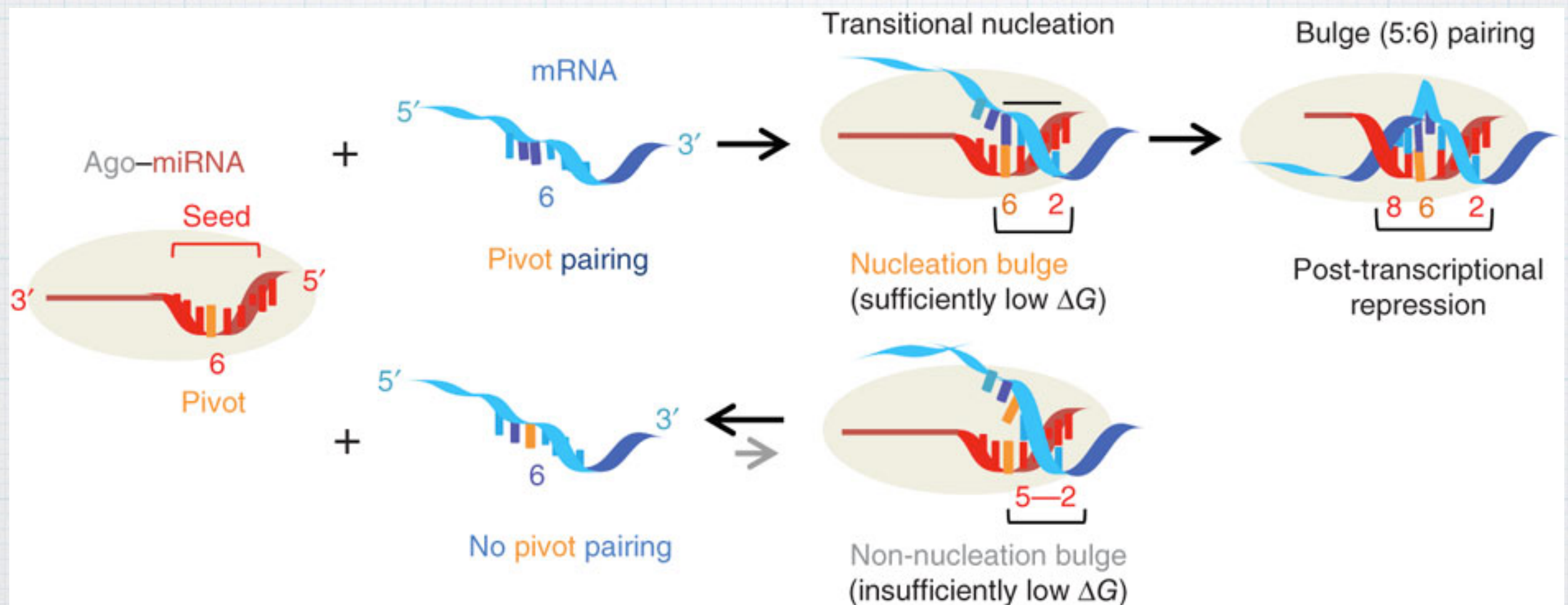
Characterization of the small RNA transcriptome using the bcbio-nextgen python framework

Lorena Pantano @lopantano
Harvard TH Chan School of Public Health

2016-07-14

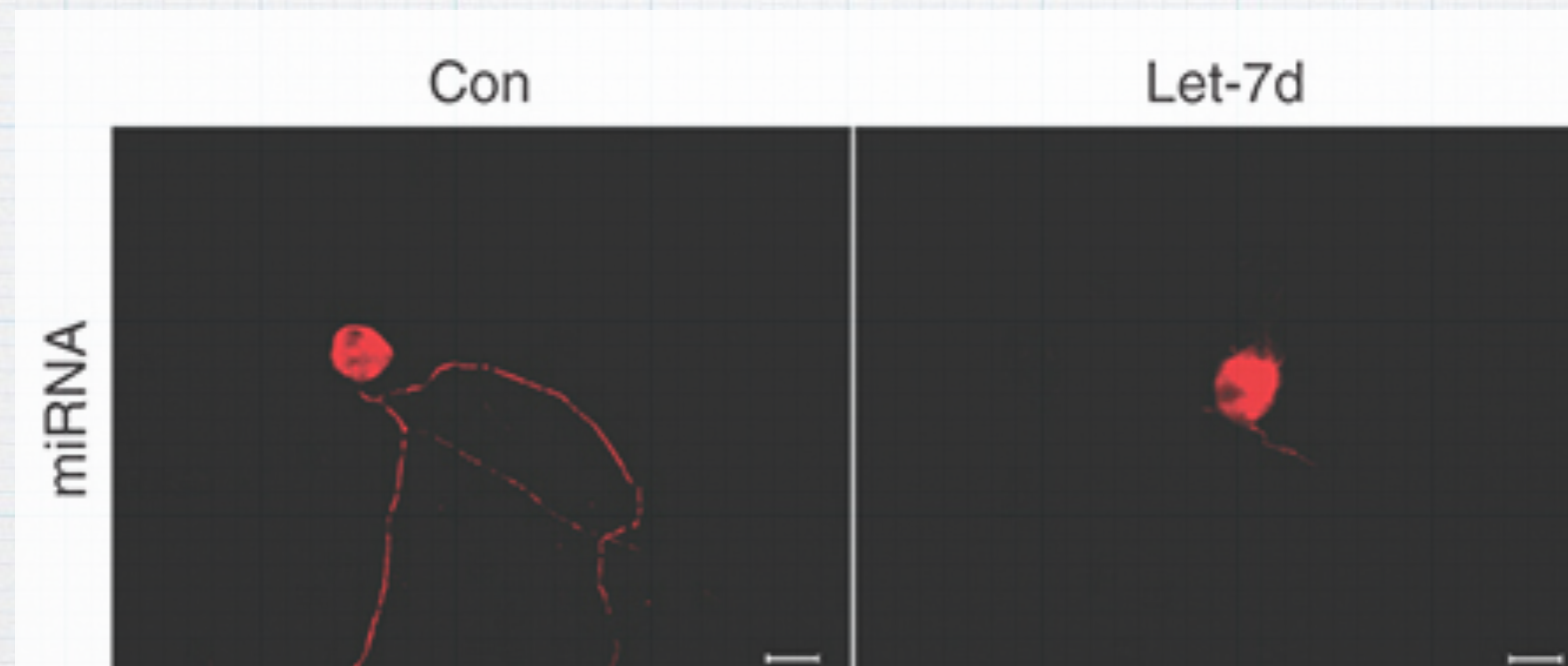
small RNA

RNA molecules of 18-36 nts
long with regulation
function



miRNA

axon outgrowth



Let-7 microRNAs Regenerate Peripheral Nerve Regeneration by Targeting Nerve Growth Factor
Shiying Li, Xinghui Wang, Yun Gu, Chu Chen, Yaxian Wang, Jie Liu, Wen Hu, Bin Yu, Yongjun Wang, Fei Ding, Yan Liu and Xiaosong Gu

isomiRs

hsa-miR-24-1-5p

hsa-miR-24-3p

```

.....GGUGCCUACUGAGCUGAUAUC.....
.....GUGCCUACUGAGCUGAUAUCAGU.....
.....GUGCCUACUGAGCUGAUAUCAG.....
.....GUGCCUACUGAGCUGAUA.....
.....UGCCUACUGAGCUGAUAUCA.....
.....UGCCUACUGAGCUGAUAUCAGU.....
.....UGCCUACUGAGCUGAUAUC.....
.....UGCCUACUGAGCUGAUA.....
.....CCUACUGAGCUGAUAUCA.....
.....CCUACUGAGCUGAUAUCAGU.....
.....CUACUGAGCUGAUAUCA.....
.....CUACUGAGCUGAUAUC.....

```

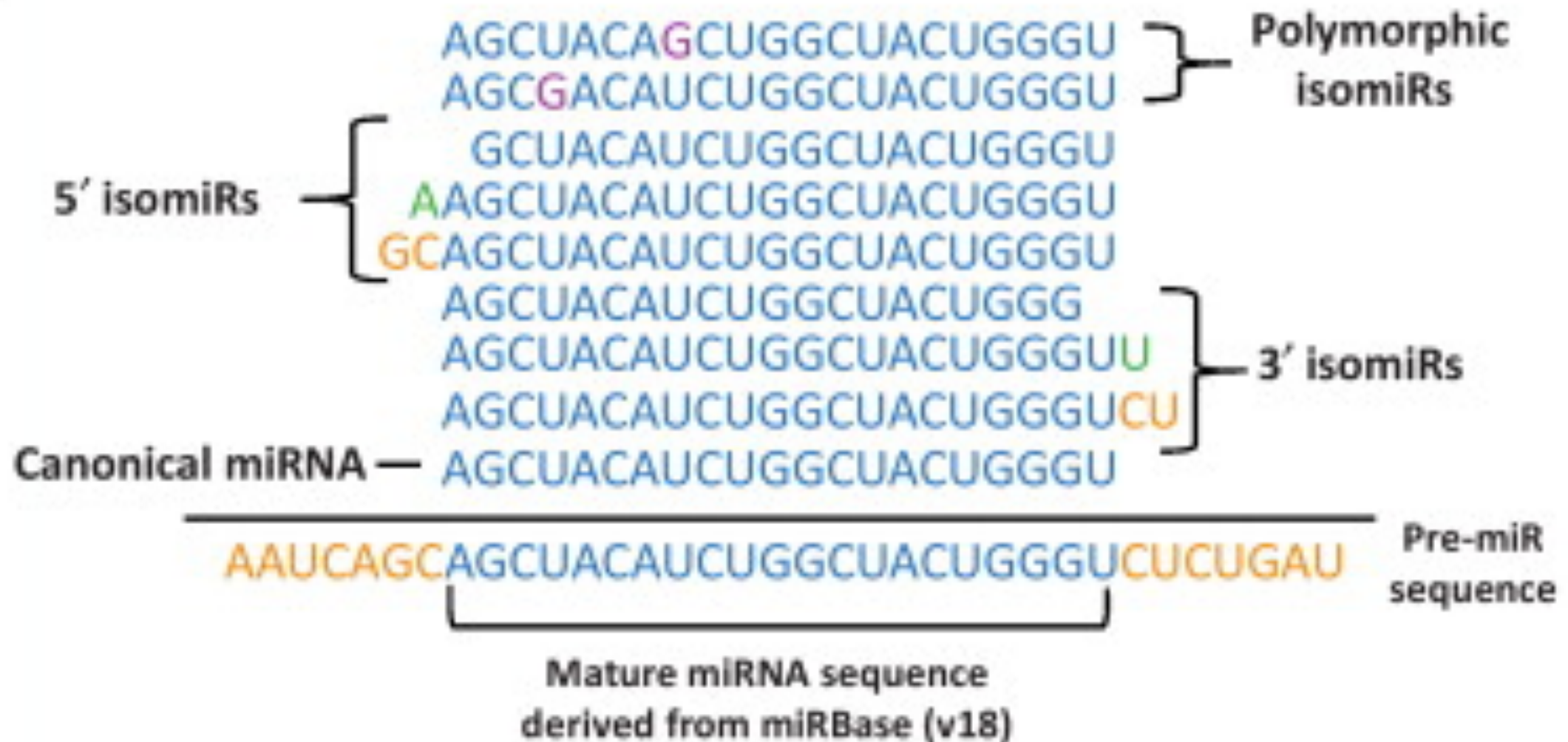
```

CUCCGGUGCCUACUGAGCUGAUAUCAGUUCUCAUUUUACACACUGGCUCAGUUCAGCAGGAACAGGAG
(((((((((.....)))))))))(((((.....))))))(((((.....))))))(((((.....))))))(-26.32)

```

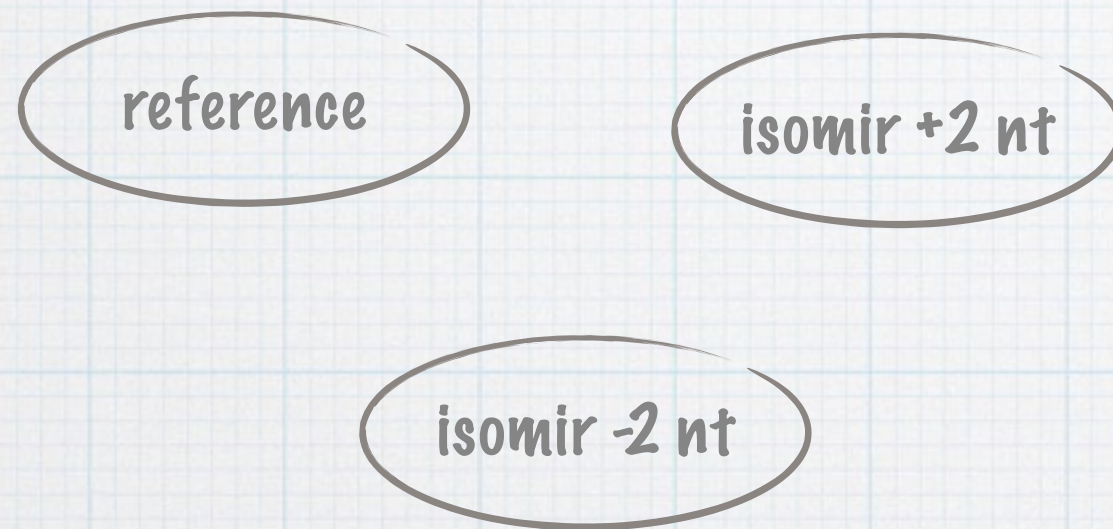
precursor

types of isomiRs

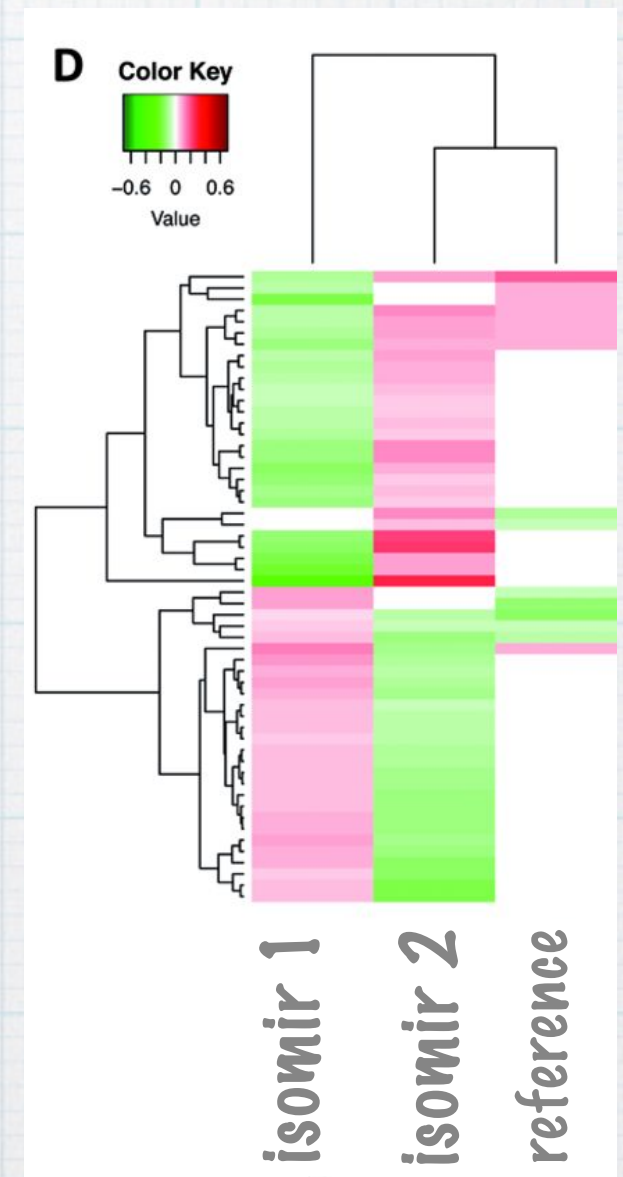


isomiRs

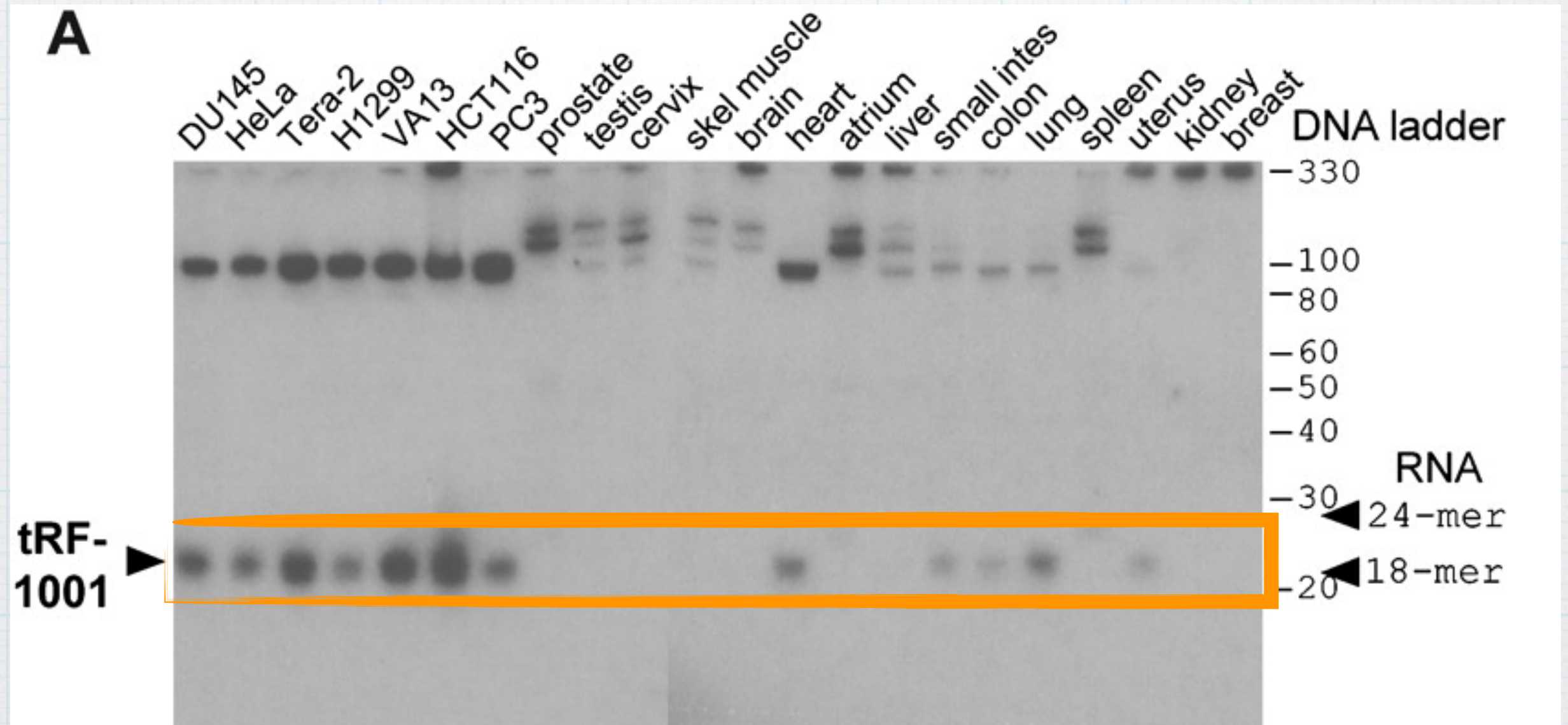
Gene expression



transfected mammary cells line
derived from metastatic site

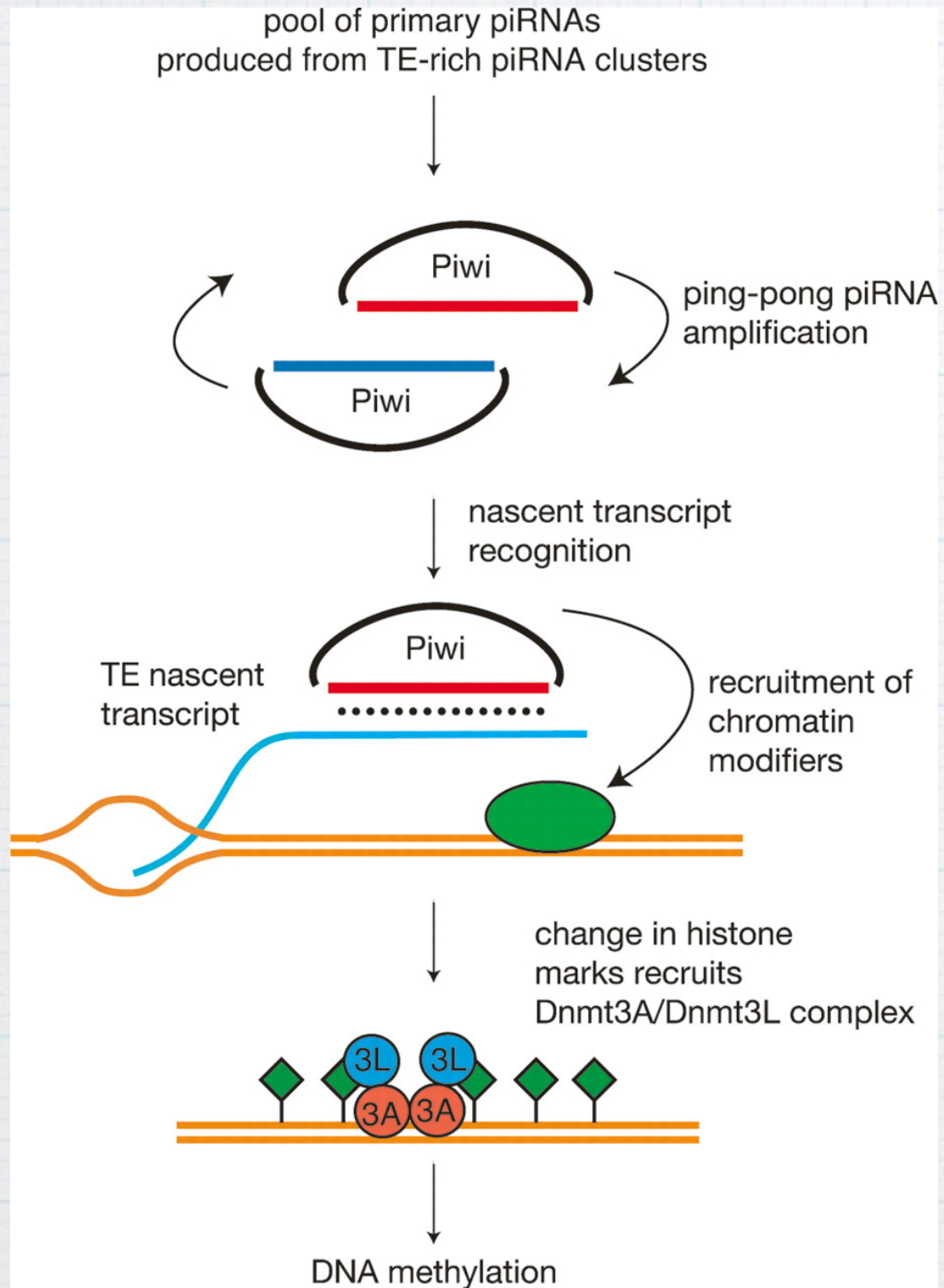


small tRNAs



Yong Sun Lee et al. Genes Dev. 2009;23:2639-2649

piRNAs

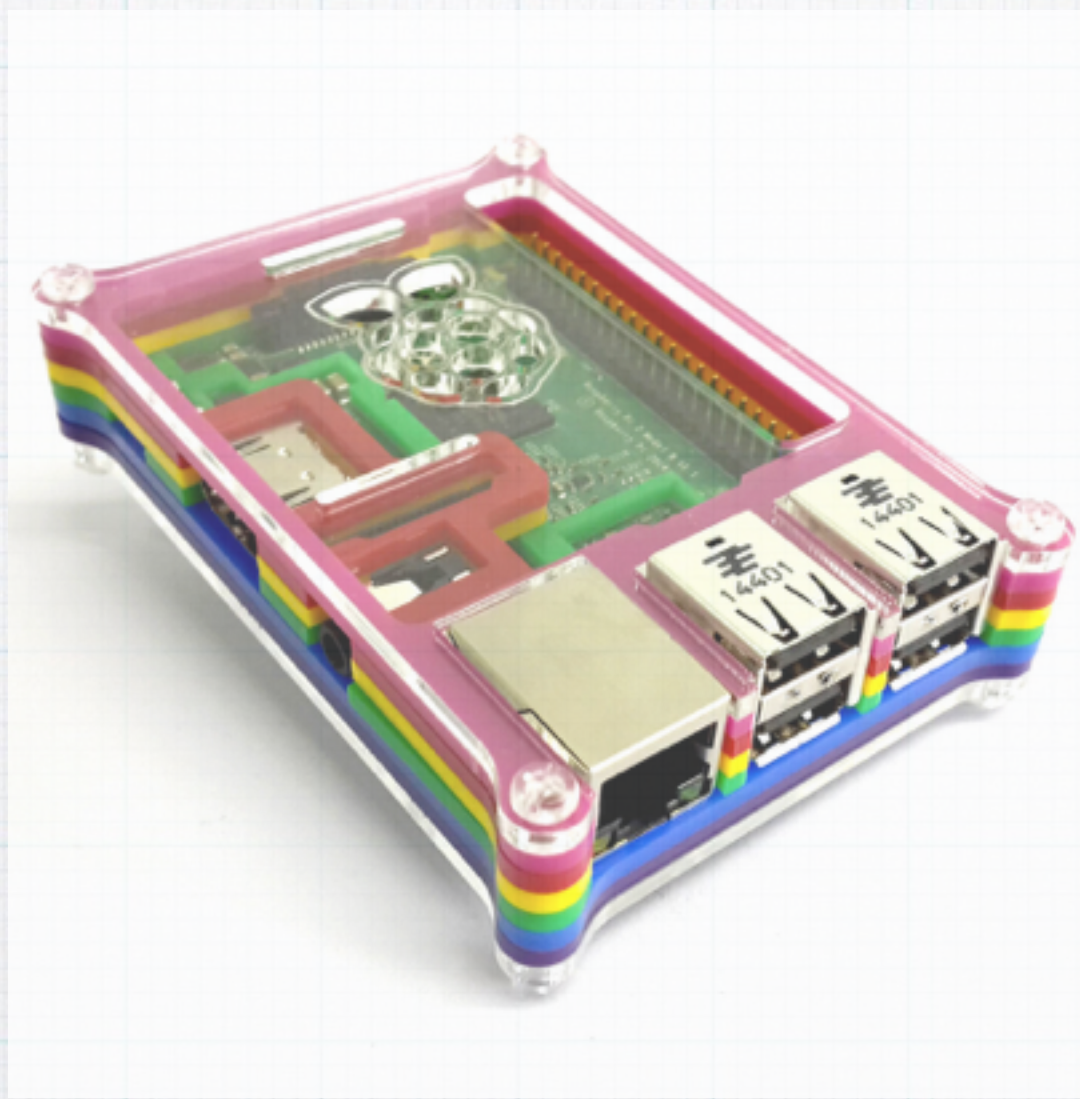


Alexei A. Aravin, and Déborah Bourc'h *Genes Dev.* 2008;22:970-975

challenges

- * isomiRs detection
- * small RNAs coming from multiple precursors over the genome (**multi-mapped reads can be 40% of the data.**)
- * differentiate degradation and functional molecules
- * non-model organism

bbio-nextgen



Variant calling, RNA-seq, small RNA-seq
over 200 peer reviewed tools **BIOCONDA[®]**

May 29, 2016 – June 29, 2016

Period: 1 month ▾

Overview

2 Active Pull Requests

74 Active Issues

2

Merged Pull Requests

0

Proposed Pull Requests

62

Closed Issues

12

New Issues

Excluding merges, **5 authors** have pushed **66 commits** to master and **66 commits** to all branches. On master, **68 files** have changed and there have been **1,085 additions** and **393 deletions**.



small RNA-seq analysis

processing & QC

cutadapt
fastqc
qualimap
multiqc

detection & annotation

miraligner
tdrmapper

de-novo

seqcluster
mirdeep2 for mirna
protac for pirna (next)

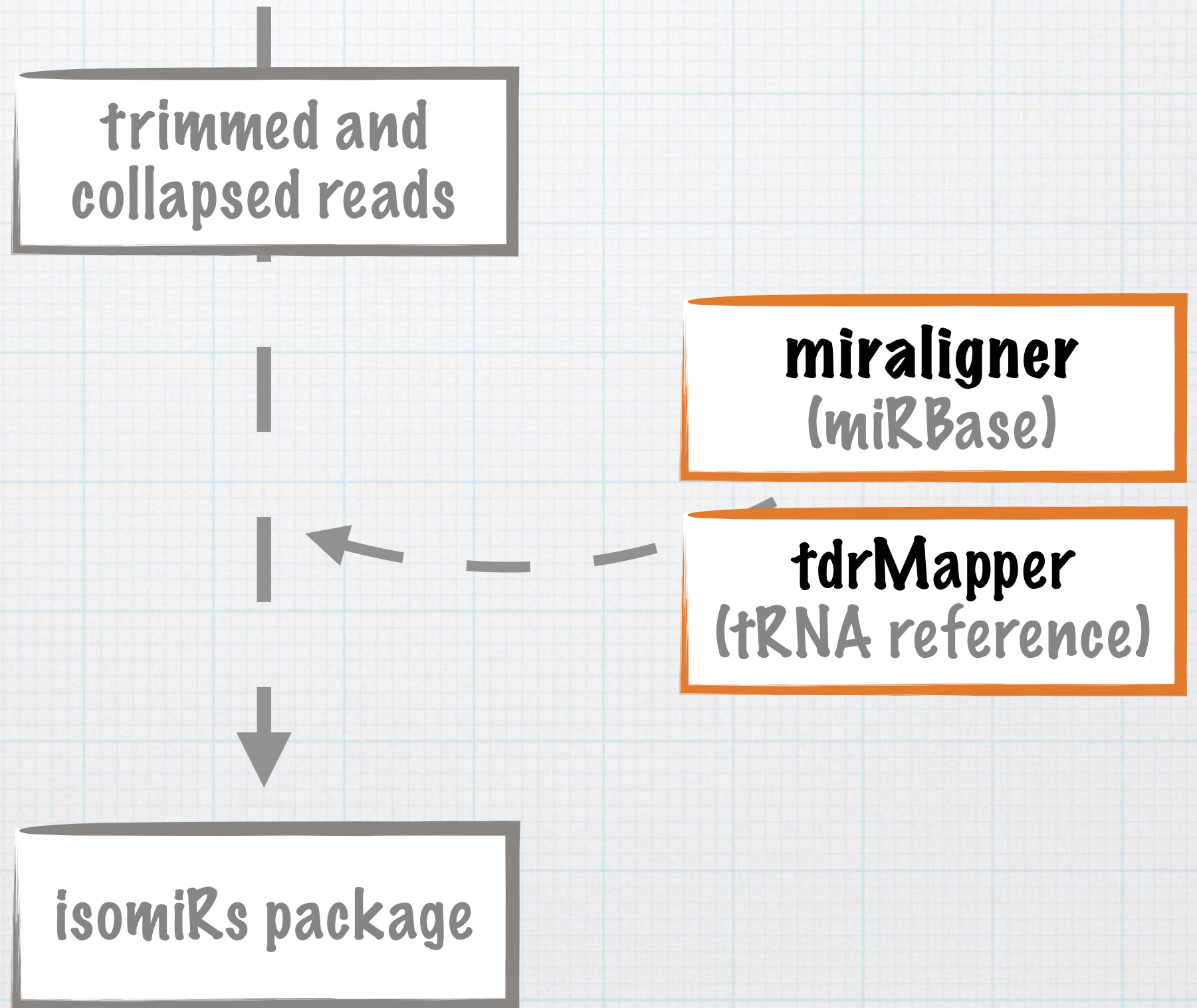
detection & annotation

trimmed and
collapsed reads

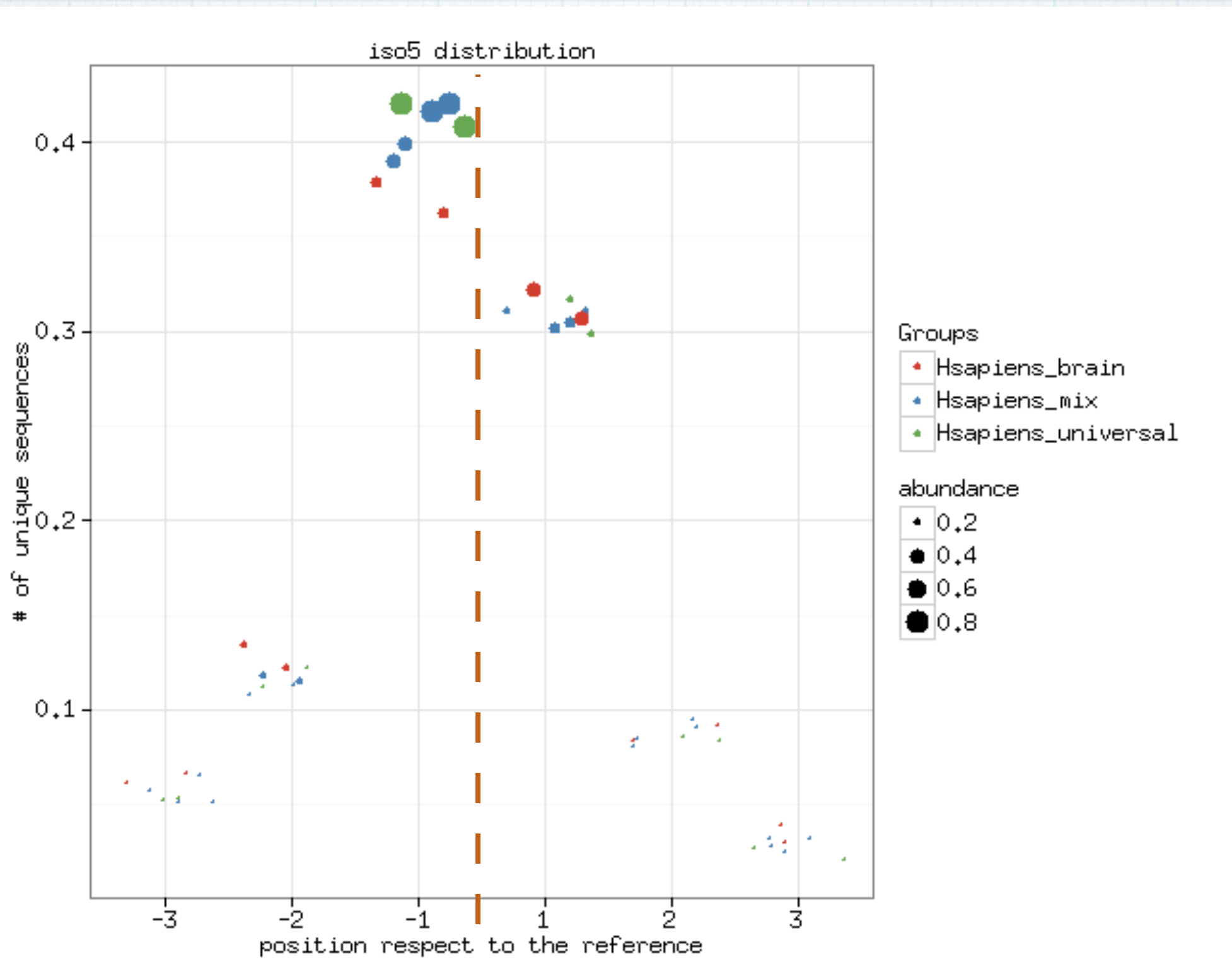
miraligner
(miRBase)

tdrMapper
(tRNA reference)

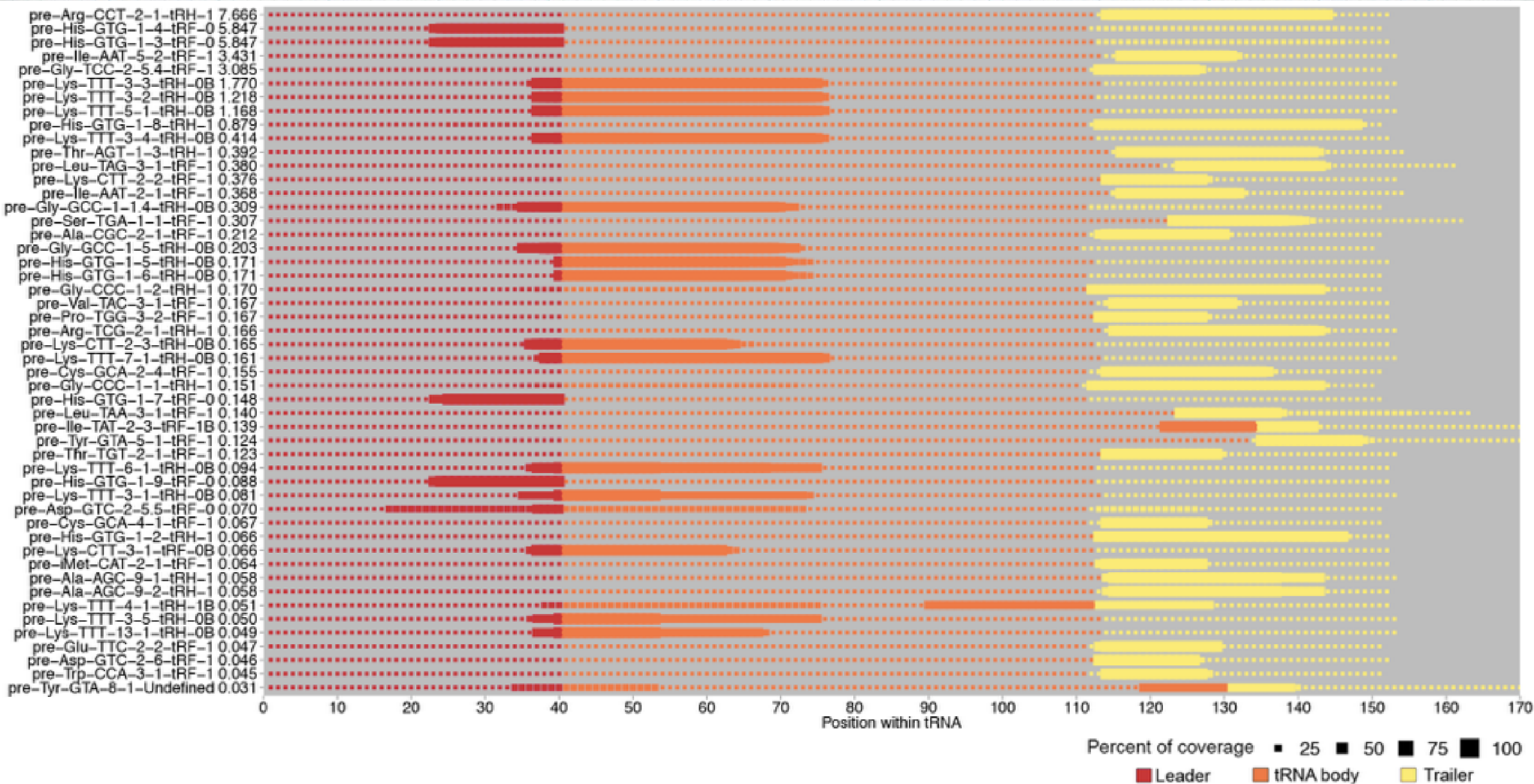
isomiRs package



isomiRs at 5' end of the miRNAs



tRNA analysis



*Pre-tRNA coverage map from NIH roadmap H1 derived mesendoderm cells, accession ID: GSM1296464

de-novo detection

trimmed and
collapsed reads

collapsing samples
into one

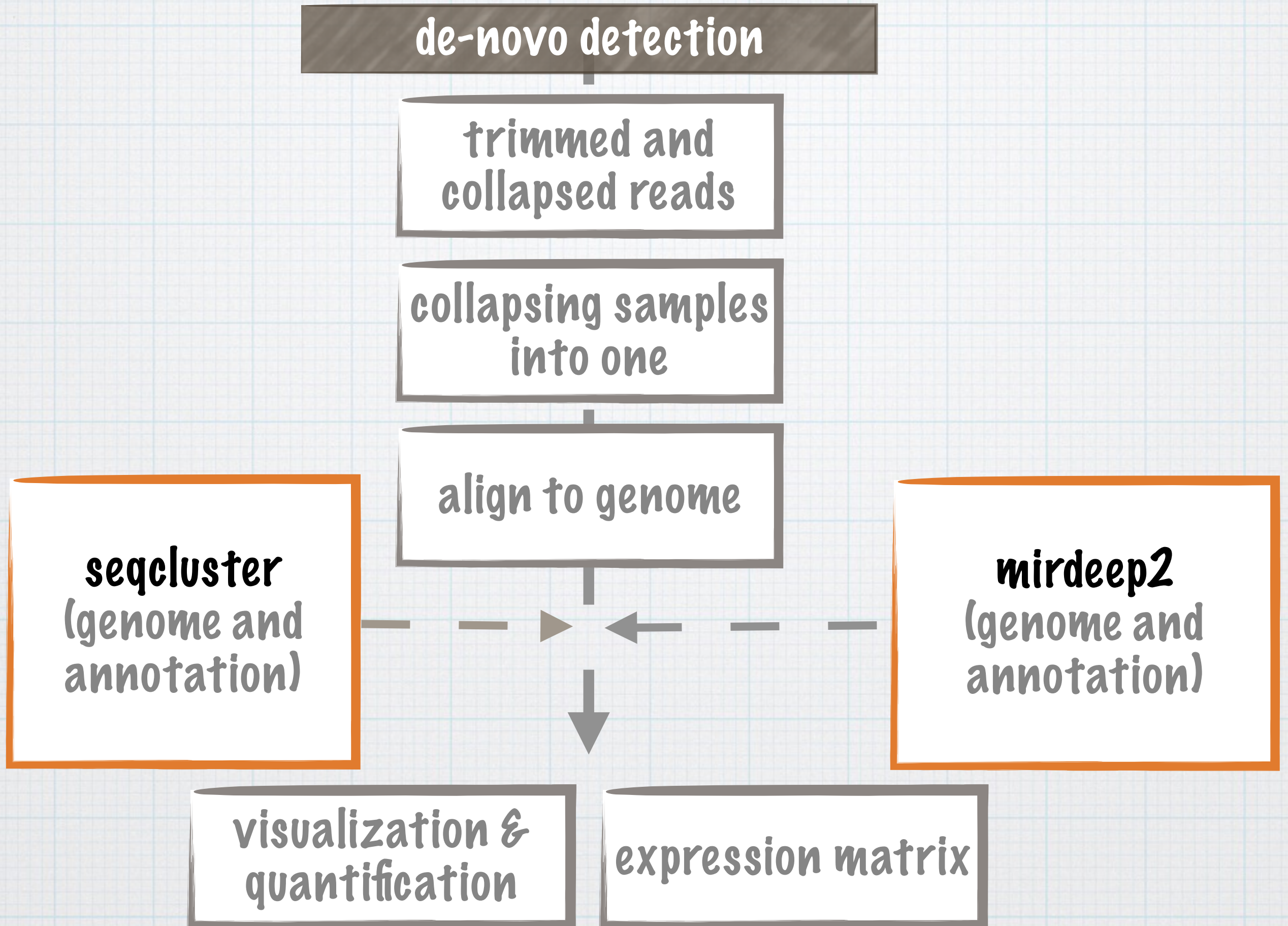
align to genome

seqcluster
(genome and
annotation)

mirdeep2
(genome and
annotation)

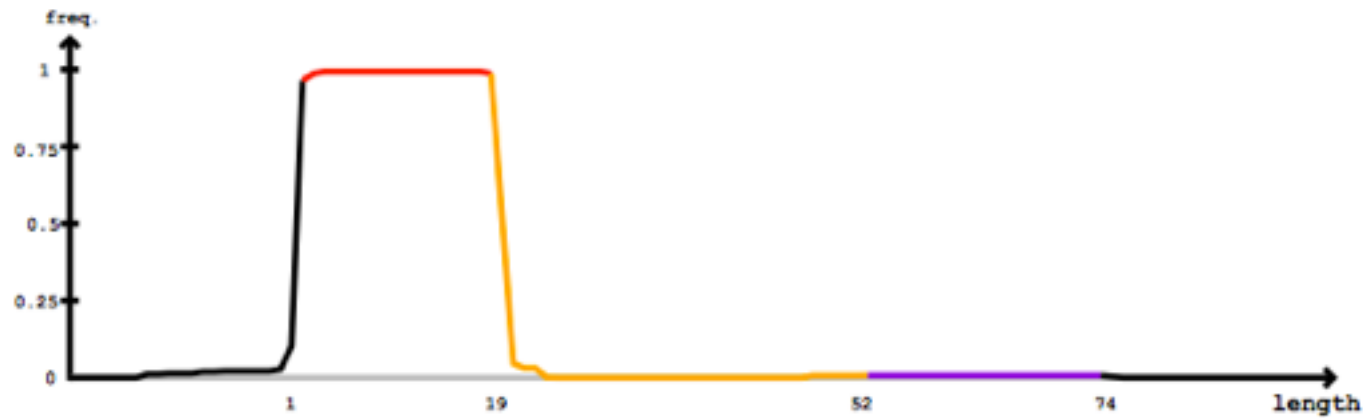
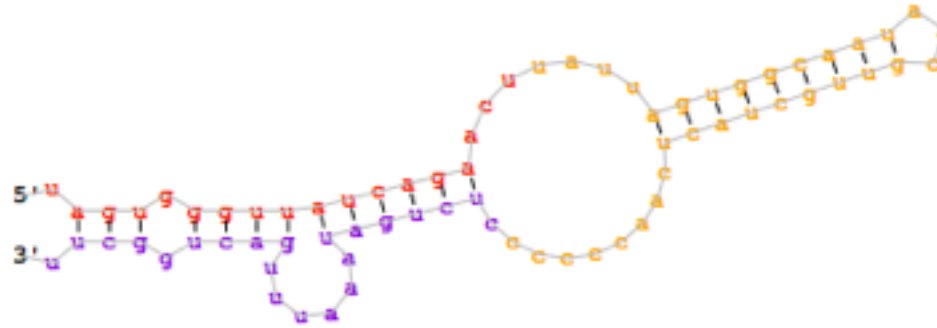
visualization &
quantification

expression matrix



miRDeep2 output

Provisional ID	: chr12_16160
Score total	: 1869.4
Score for star read(s)	: 3.9
Score for read counts	: 1866.6
Score for mfe	: -1
Score for randfold	:
Score for cons. seed	:
Total read count	: 3673
Mature read count	: 3670
Loop read count	: 0
Star read count	: 3

[illegible]

seqcluster



The diagram illustrates a 'meta-cluster' containing two distinct clusters of multi-mapped reads. Each cluster is represented by a group of four blue horizontal lines (multi-mapped reads) and one grey horizontal line (single-mapped read). The top cluster is labeled 'cluster at position 1' and the bottom cluster is labeled 'cluster at position 2'. The entire set of reads is enclosed in a black rectangular frame.

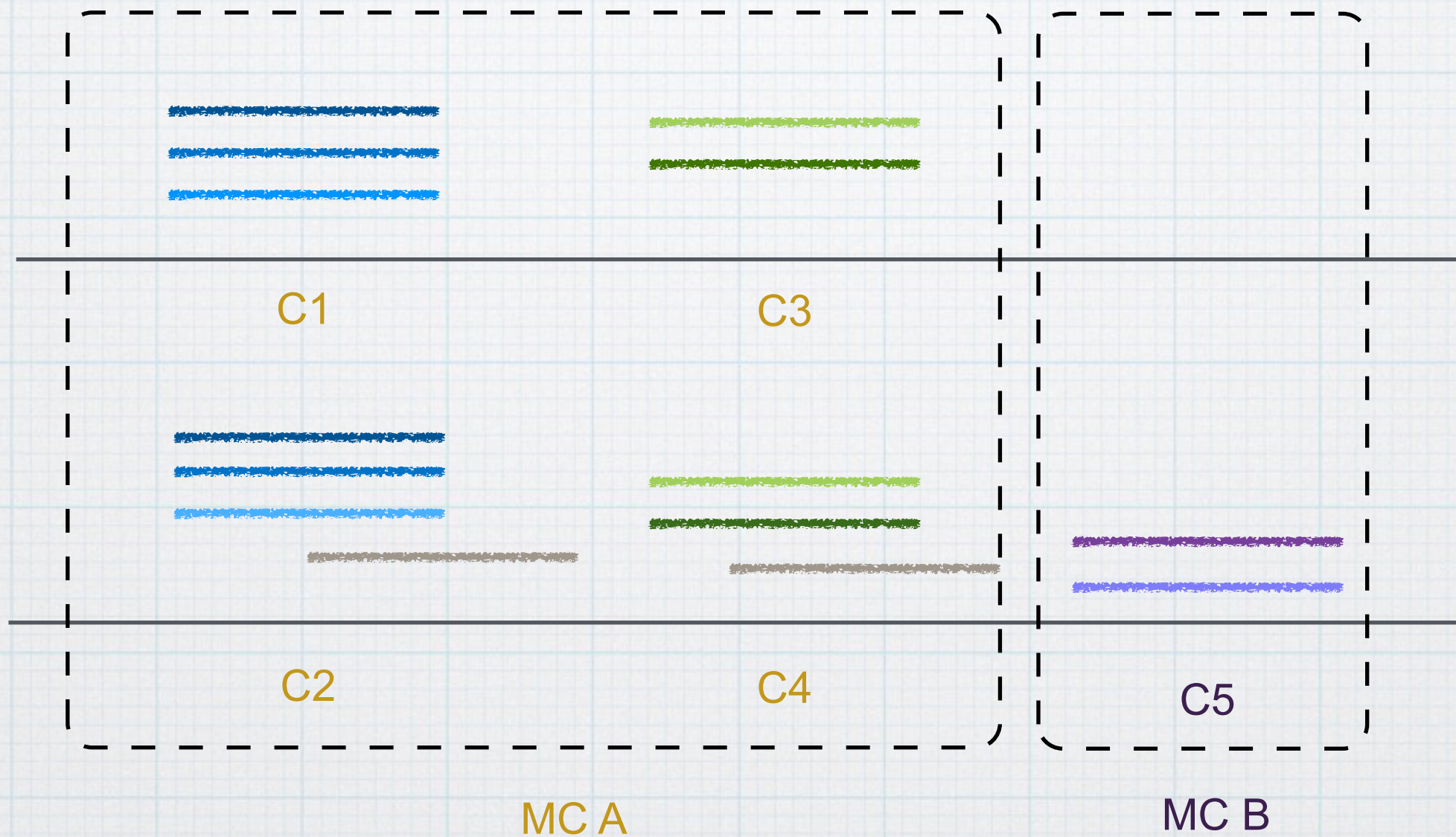
cluster at position 1

cluster at position 2

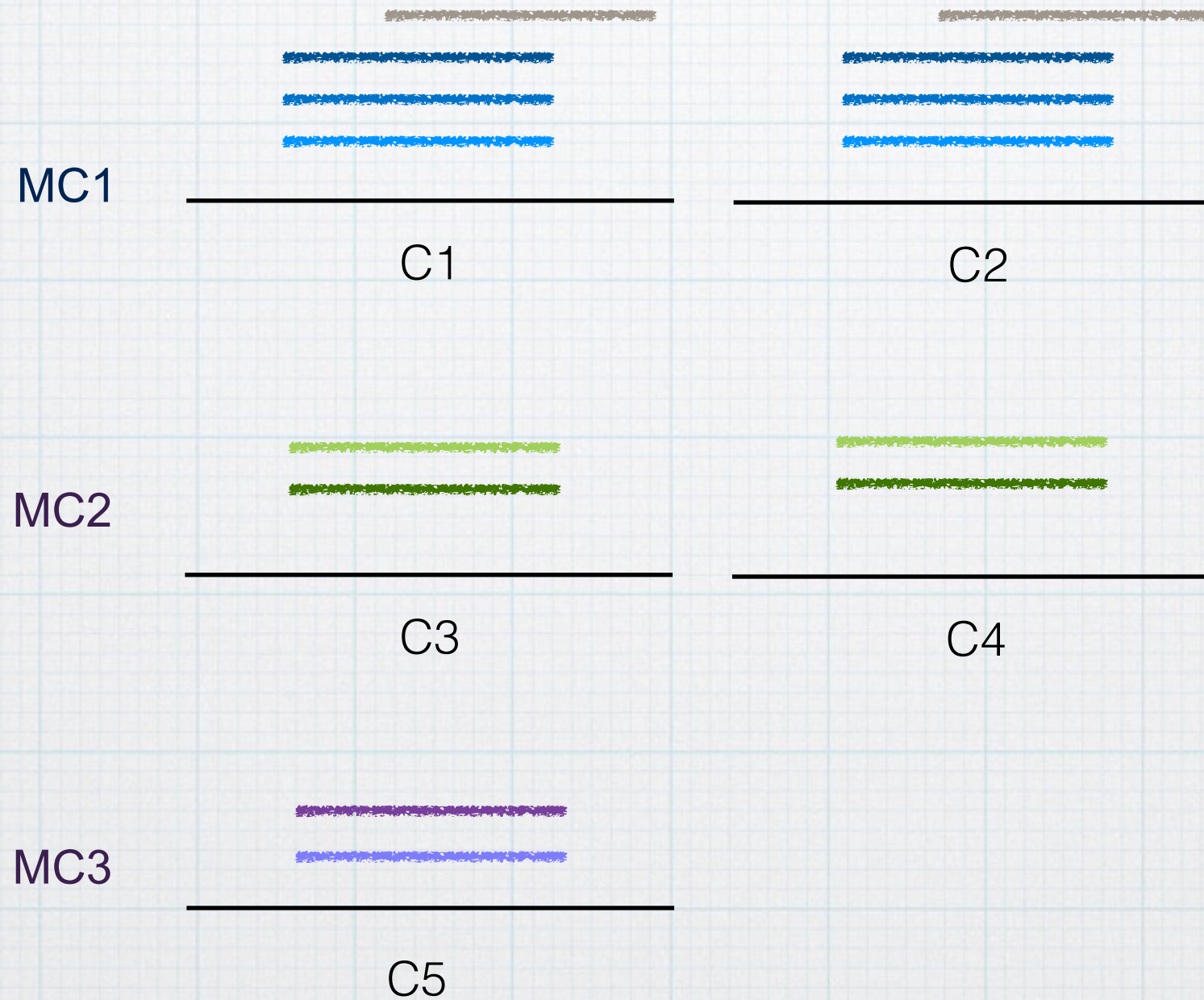
meta-cluster

seqcluster deals with multi-mapped reads

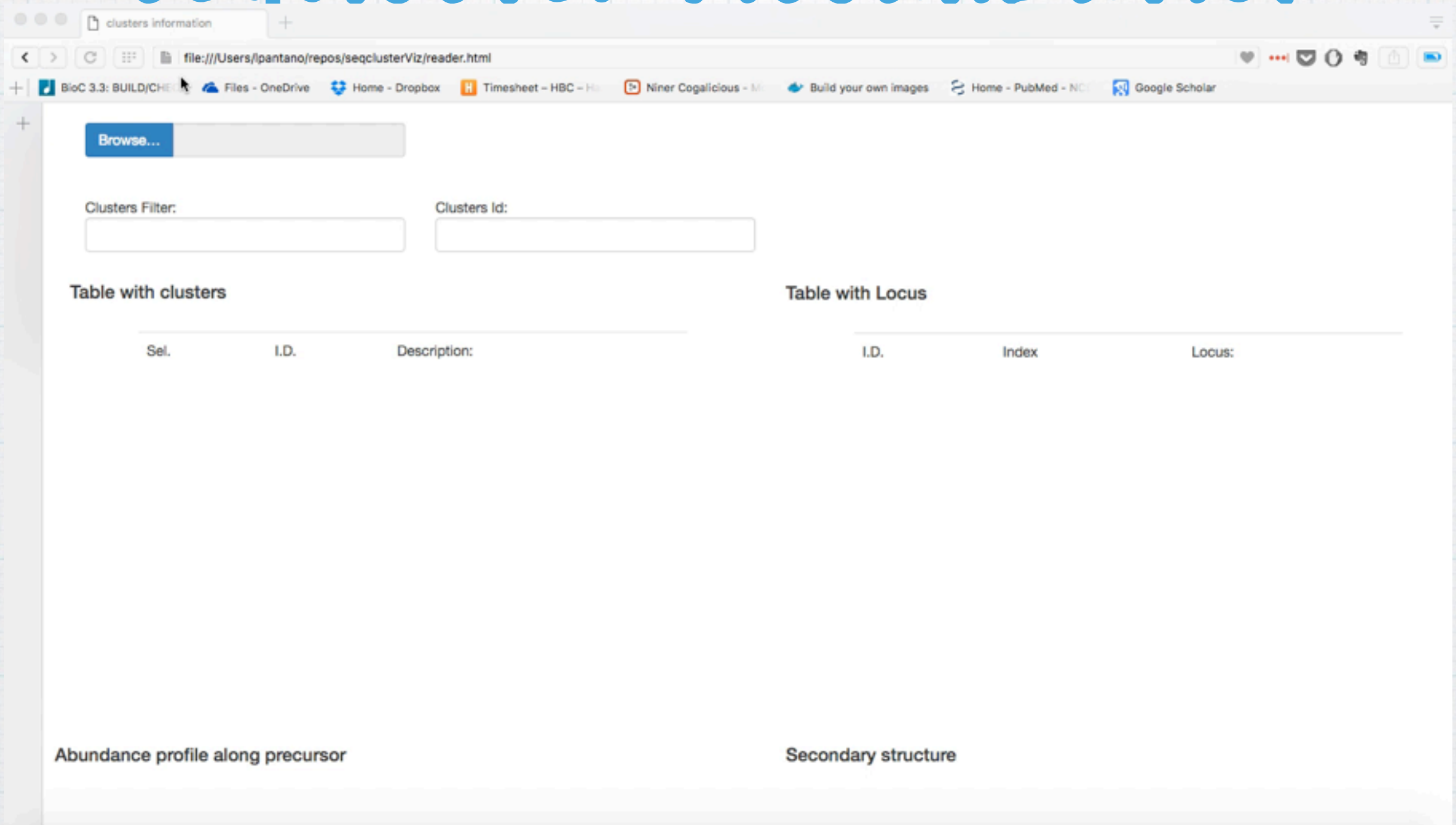
Step 1: clustering



Step 2: cleaning



seqcluster visualization



<https://github.com/lpantano/seqclusterViz>

MultiQC



Phil Ewels
ewels

Bioinformatician working with next generation sequencing data.

Science for Life Laboratory
 Stockholm, Sweden
 phil.ewels@scilifelab.se
 <http://phil.ewels.co.uk>
 Joined on Nov 3, 2010

48 Followers
21 Starred
23 Following

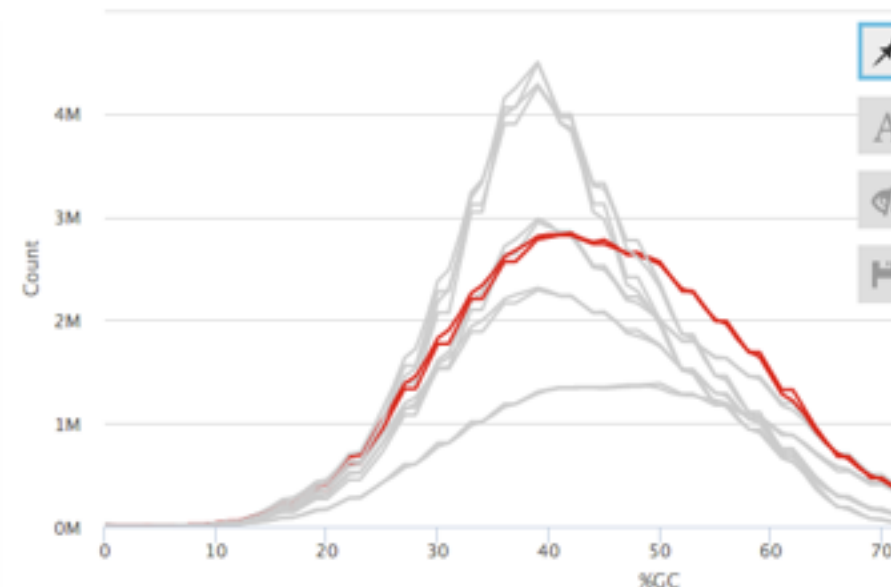
STAR: % Uniquely mapped reads		
signed	% Aligned	M Aligned
0.9	81.1%	
1.5	79.1%	
1.9	70.2%	
0.9	63.2%	
0.7	61.8%	
0.6	50.6%	

Sequence GC Content

6 5 1

age GC content of reads. Normal random library typically have a roughly normal di
GC help.

Per Sequence GC Content



MultiQC Toolbox

Highlight Samples

510

Regex mode ☐

514

miRQC project

Evaluation of quantitative miRNA expression platforms in the microRNA quality control (miRQC) study

Pieter Mestdagh, Nicole Hartmann, Lukas Baeriswyl, Ditte Andreassen, Nathalie Bernard, Caifu Chen, David Cheo, Petula D'Andrade, Mike DeMayo, Lucas Dennis, Stefaan Derveaux, Yun Feng, Stephanie Fulmer-Smentek, Bernhard Gerstmayer, Julia Gouffon, Chris Grimley, Eric Lader, Kathy Y Lee, Shujun Luo, Peter Mouritzen, Aishwarya Narayanan, Sunali Patel, Sabine Peiffer, Silvia Rüberg, Gary Schroth  *et al.*

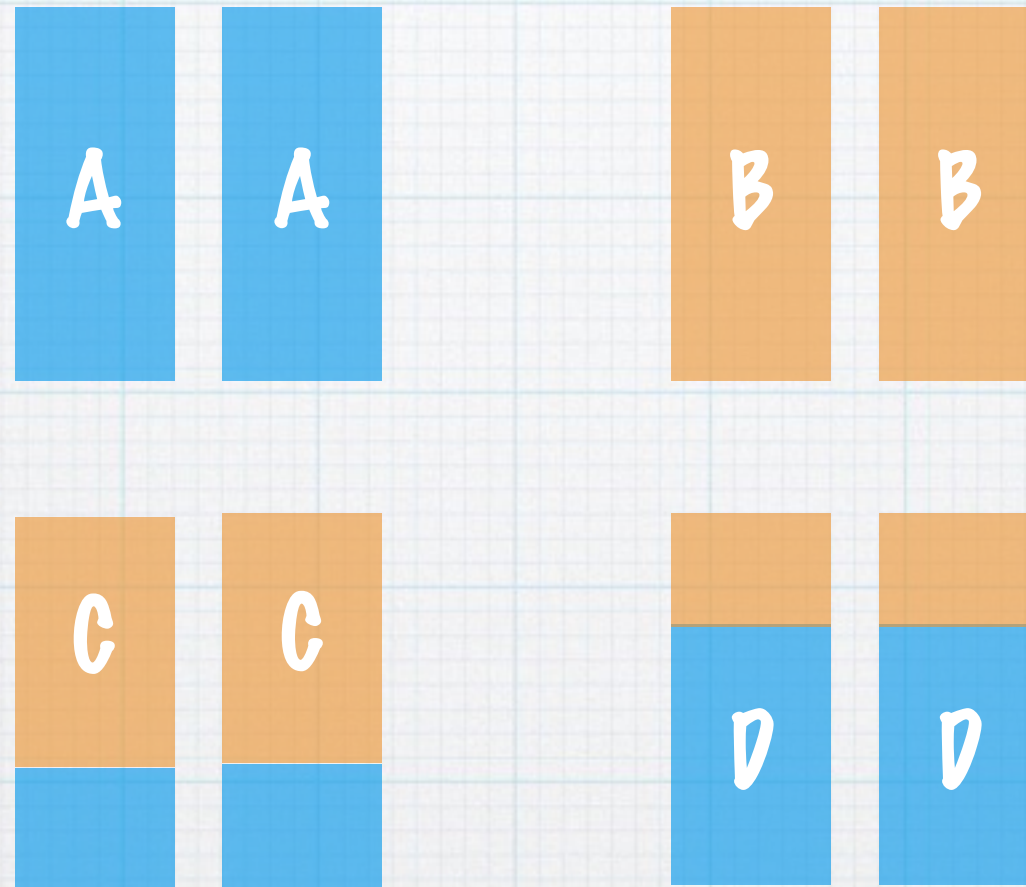
Affiliations | **Contributions** | **Corresponding author**

Nature Methods **11**, 809–815 (2014) | doi:10.1038/nmeth.3014

Received 27 February 2014 | Accepted 22 May 2014 | Published online 29 June 2014

| Corrected online **30 July 2014**

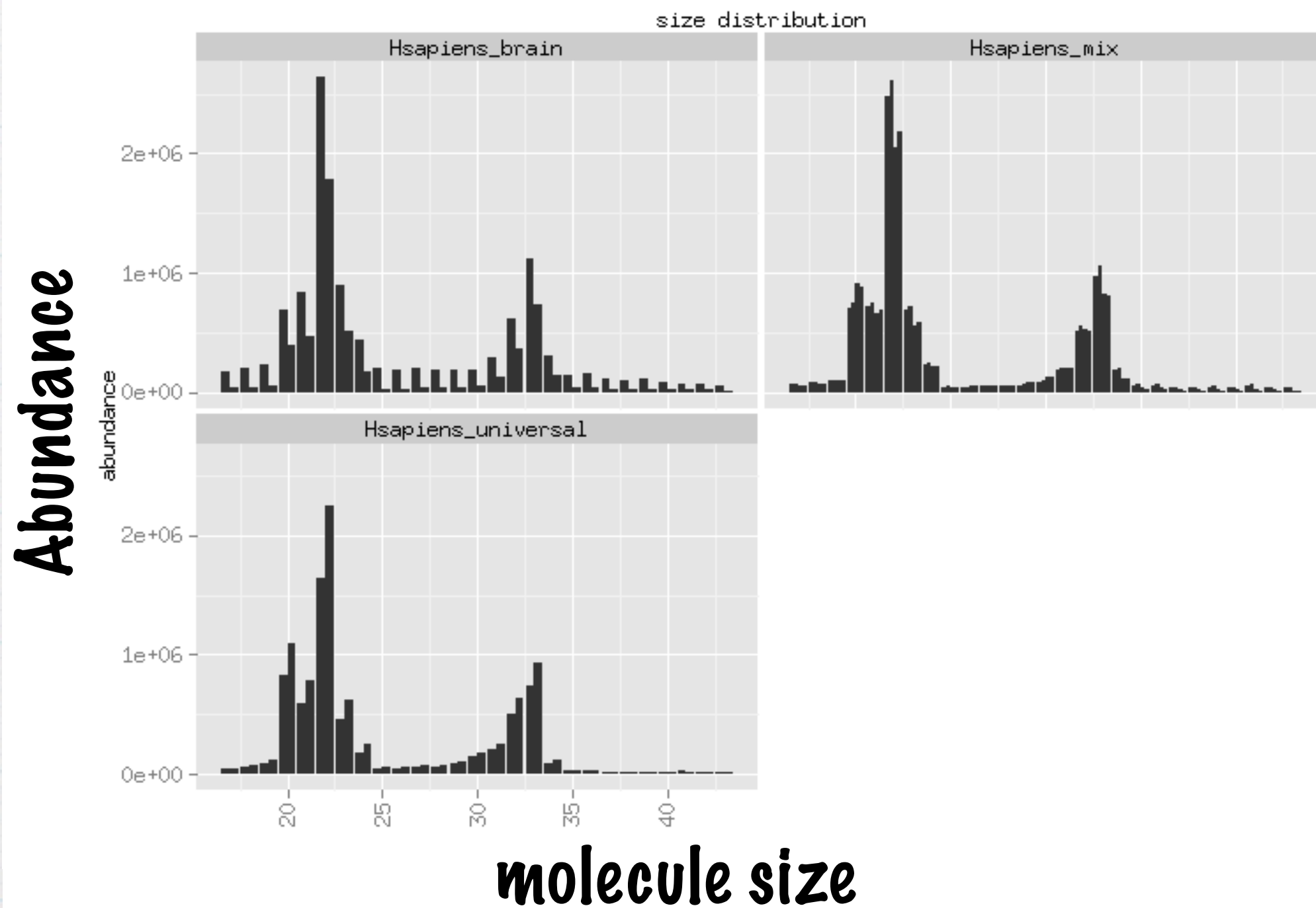
Quality Control samples



For each molecule:

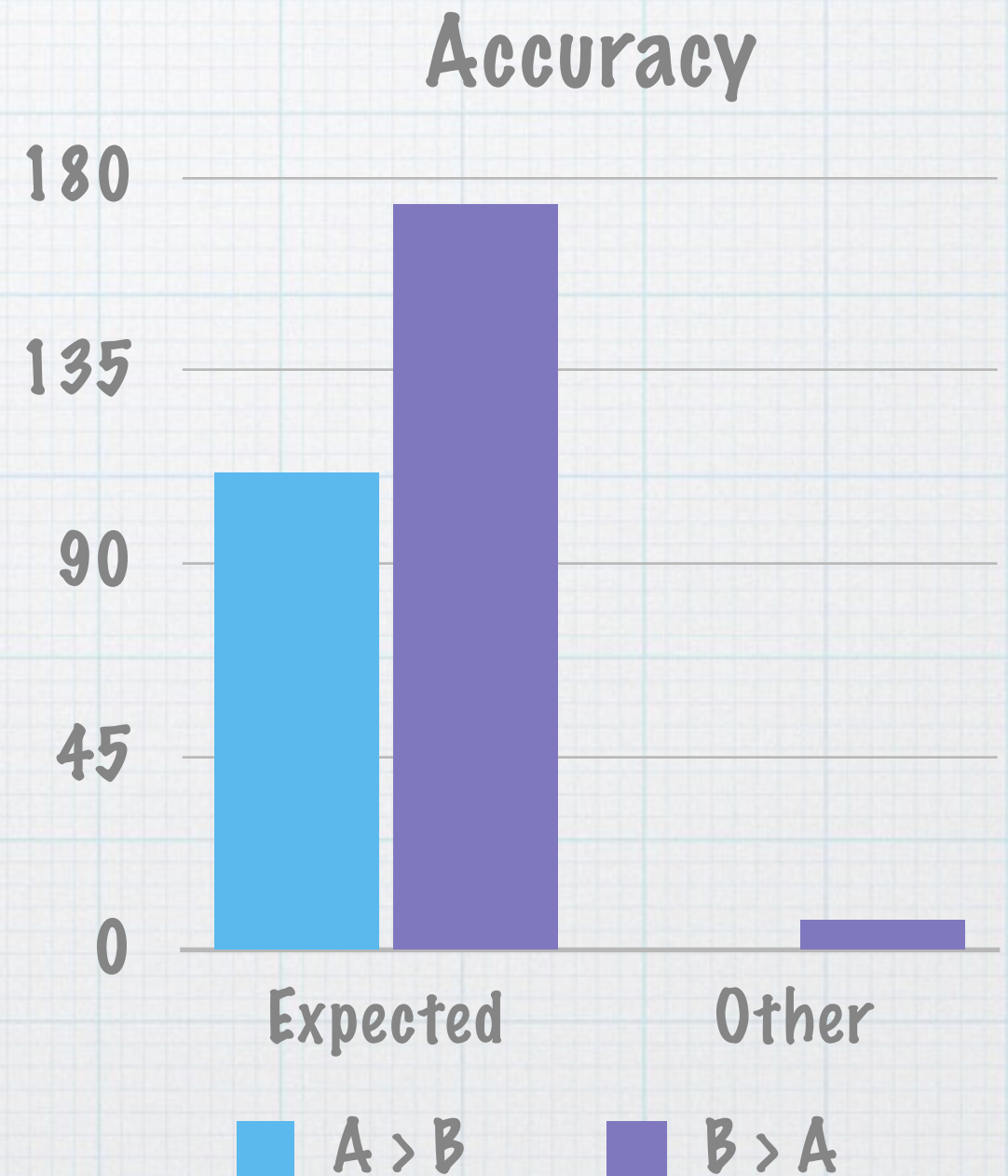
- * If $A > B$ then $A > D > C > B$
- * If $B > A$ then $A < D < C < B$

Good samples



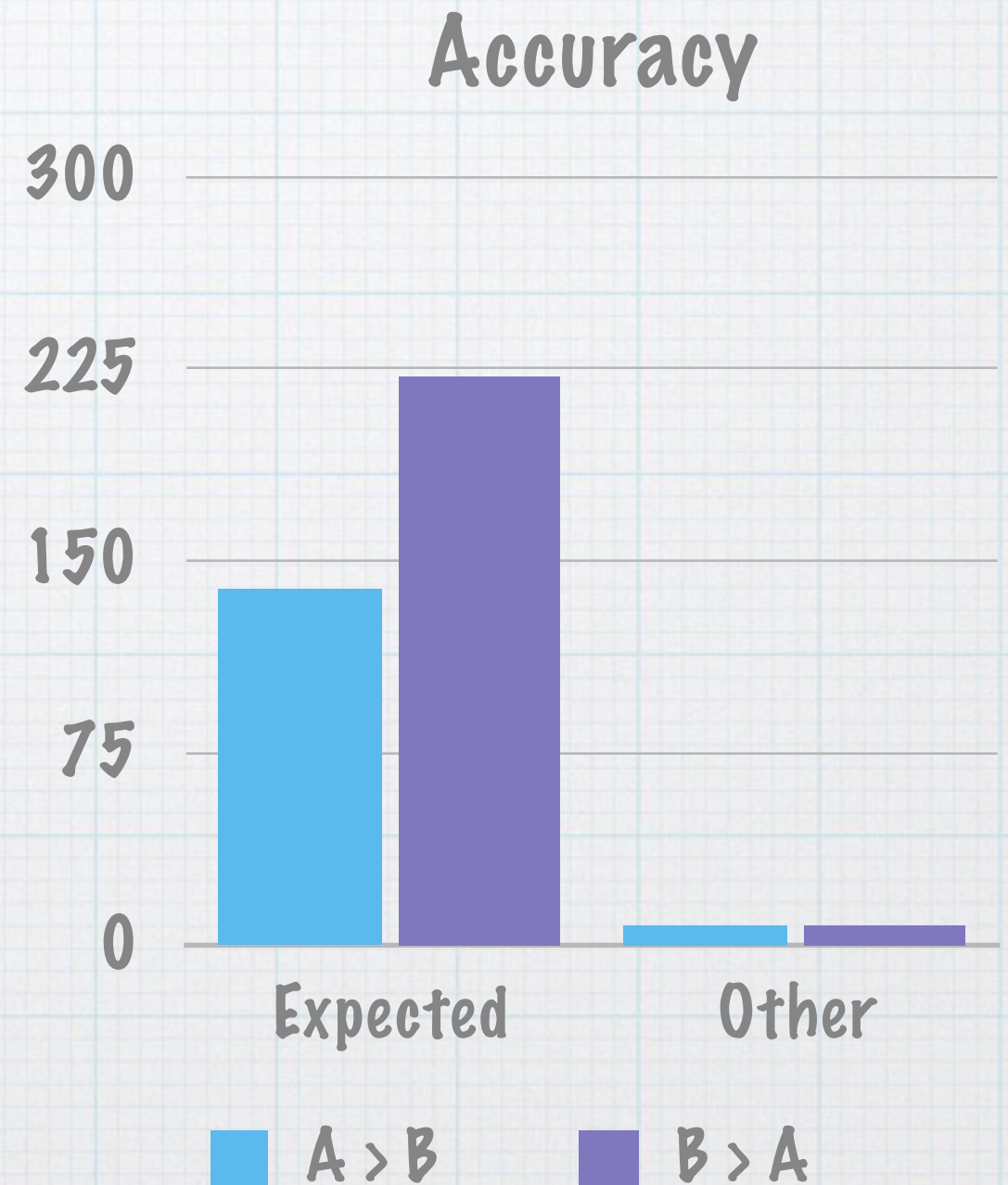
miRNA quantification

miRNAs > 5 counts in average
upper quantile normalization

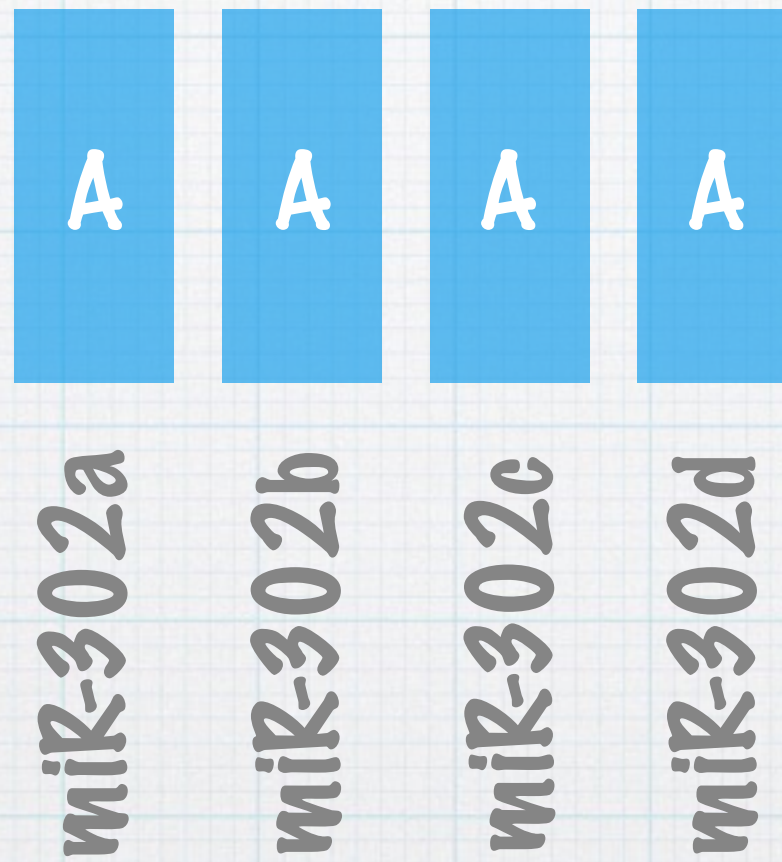


clusters quantification

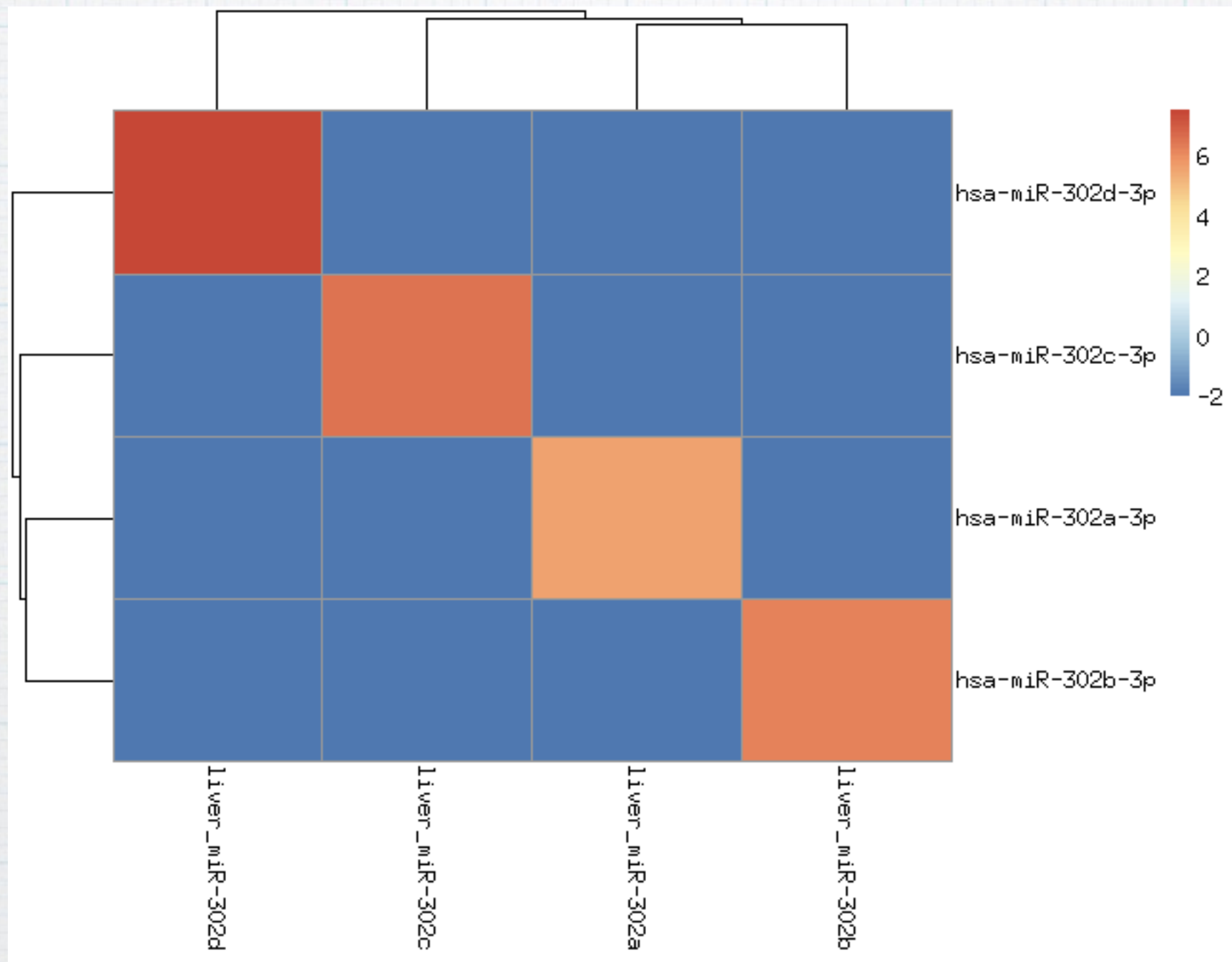
expression > 5 counts in average
upper quantile normalization



Positive controls



Specificity

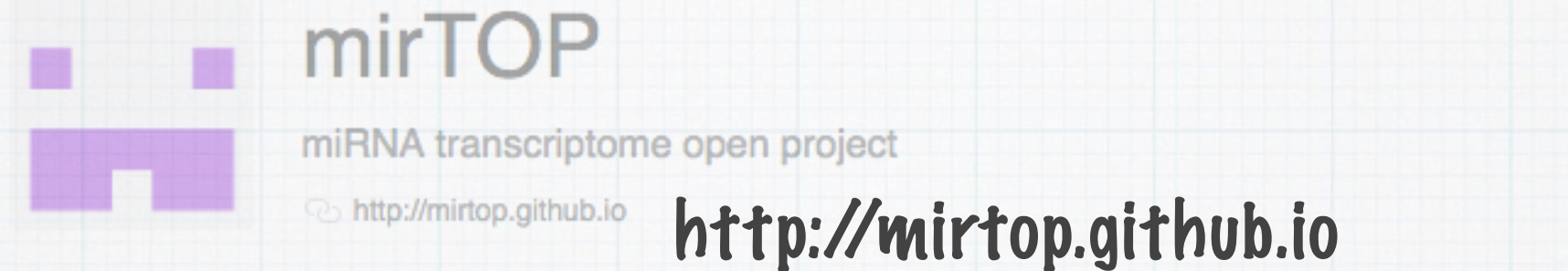


Resources

	Time (h)
organize	0:01
adapter	0:27
alignment	0:26
annotation	3:43
cluster + mirdeep2	4:15
qc	0:04

The time for 8 samples with 6 millions reads each was 8 hours and 57 minutes.

open project for small RNA annotation and analysis



<http://mirtop.github.io>

standard formats
naming rules

best-practices

miRNAs, tRNAs ...

thanks

- * Harvard T.H. Chan School of Public Health**

- * Research Computing at Harvard Medical School: Chris Botka, Director of Research Computing and all the people in the team.**

- * Special thanks to the authors of those papers to make data available.**