



small RNAseq data analysis miRNA detection

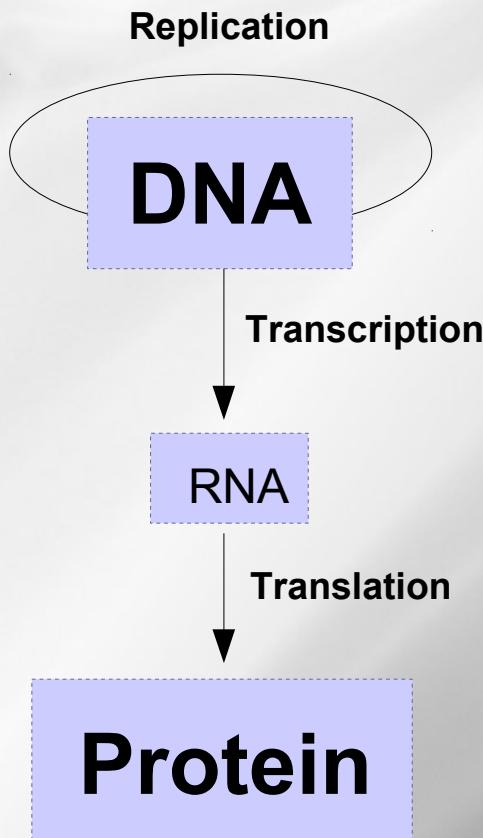
P. Bardou, W. Carré, C. Gaspin, S. Maman, J. Mariette & O. Rué

Introduction ncRNA

Central dogma of molecular biology

- **Evolution of the dogma : 1950-1970**

DNA structure discovery.

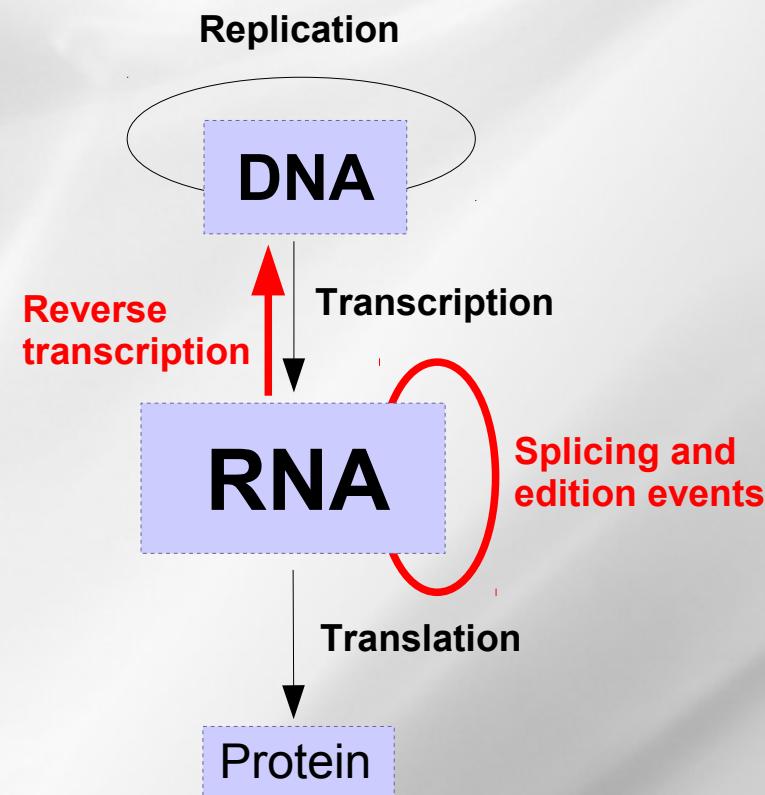


One gene = one function

Central dogma of molecular biology

- Evolution of the dogma : 1970-1980

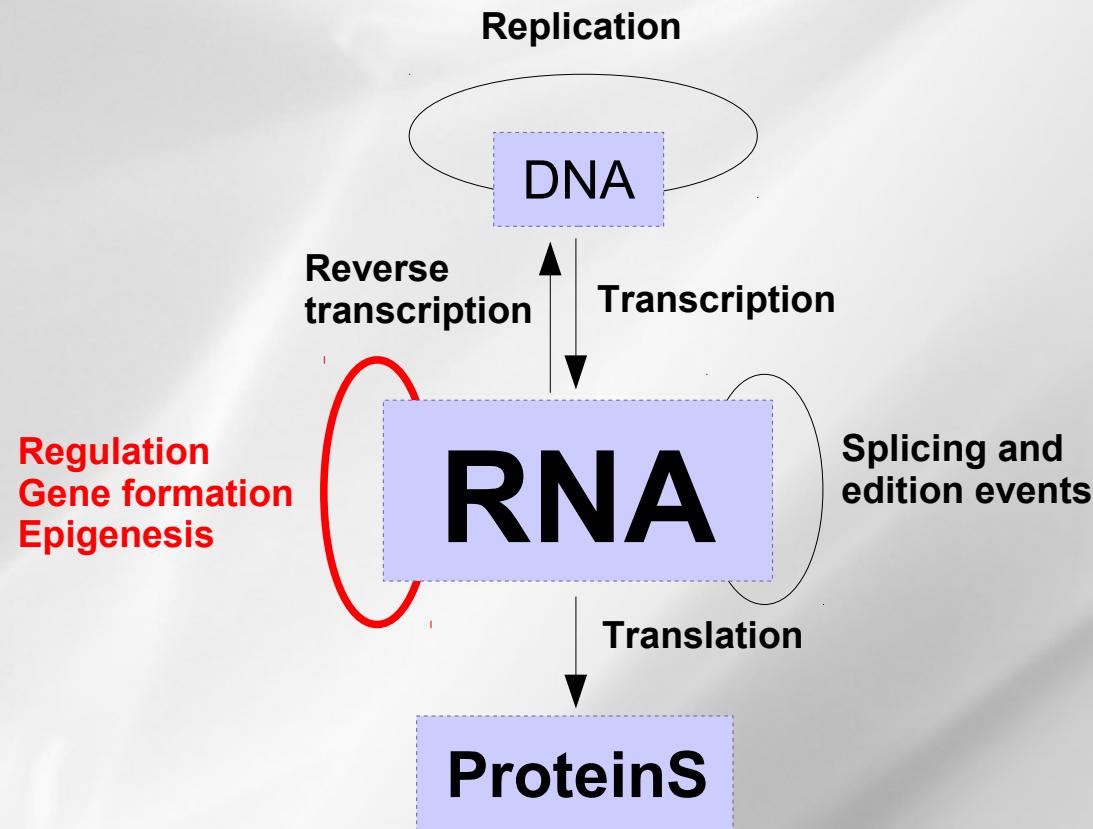
Genome analysis



Central dogma of molecular biology

- Evolution of the dogma : aujourd'hui

Genome analysis + Sequencing

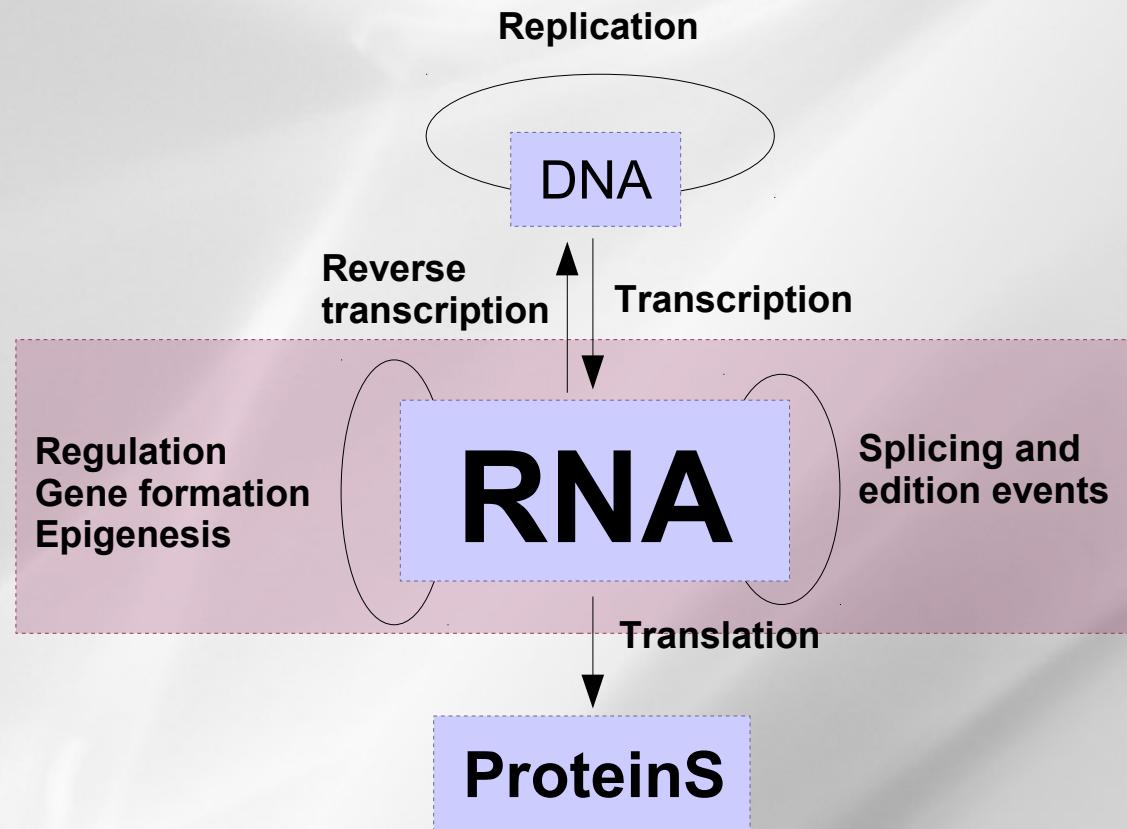


Many genes = one fonctionnel complex

Central dogma of molecular biology

- Evolution of the dogma : **aujourd'hui**

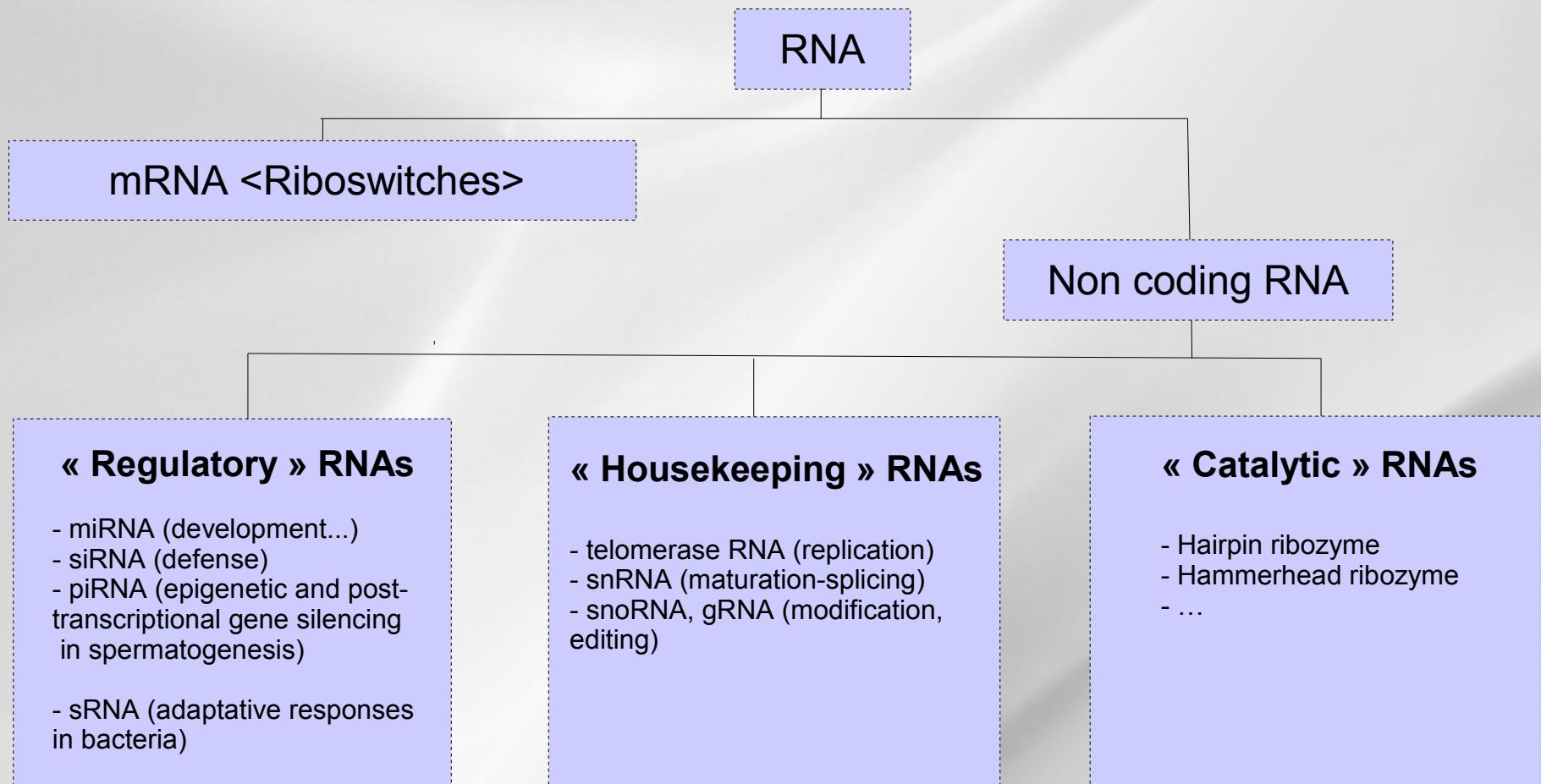
Genome analysis + Sequencing



Many genes = one fonctionnel complex

The RNA world

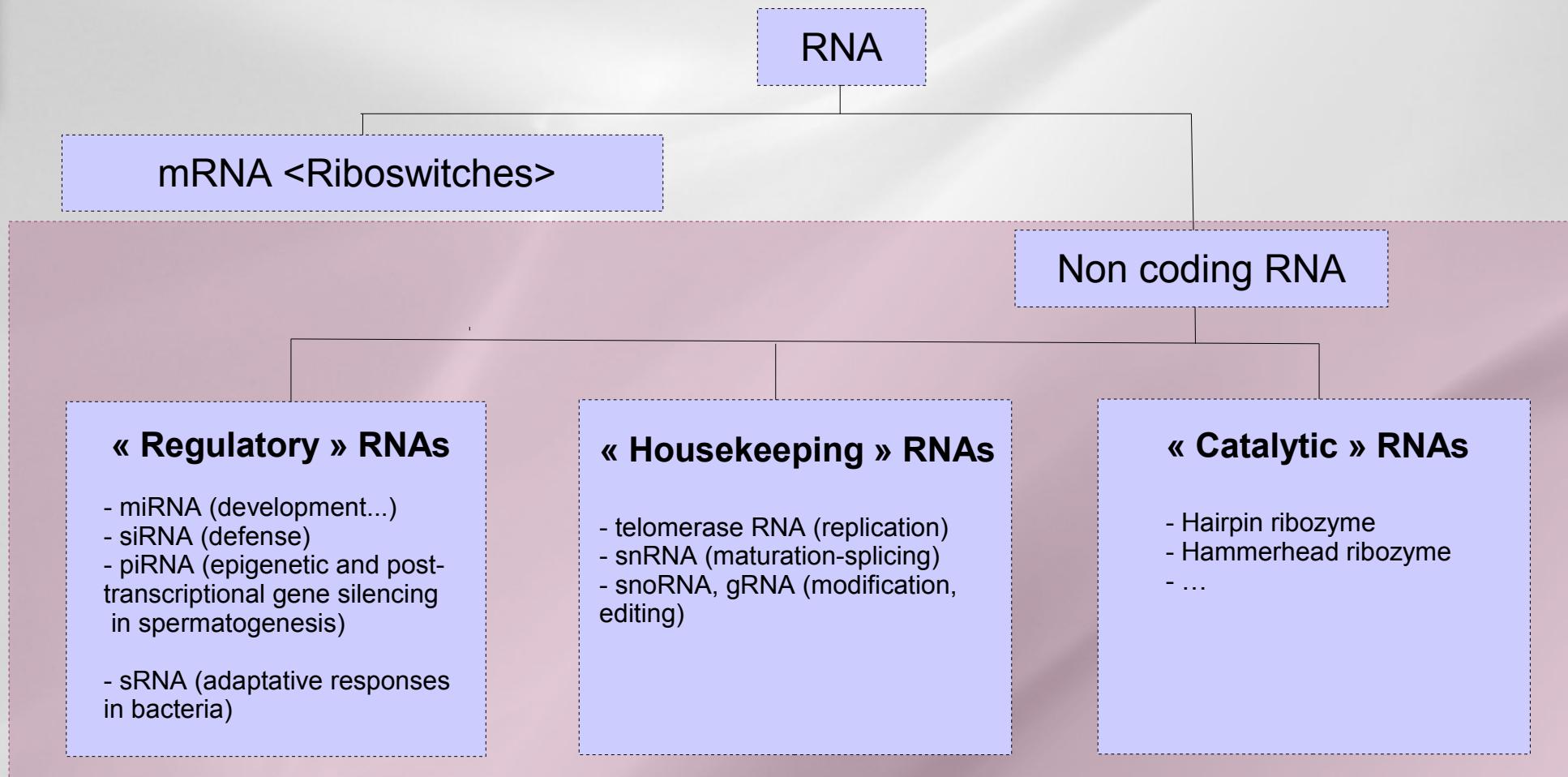
- An expending universe of RNA



→ Multiple roles of RNA in genes regulation

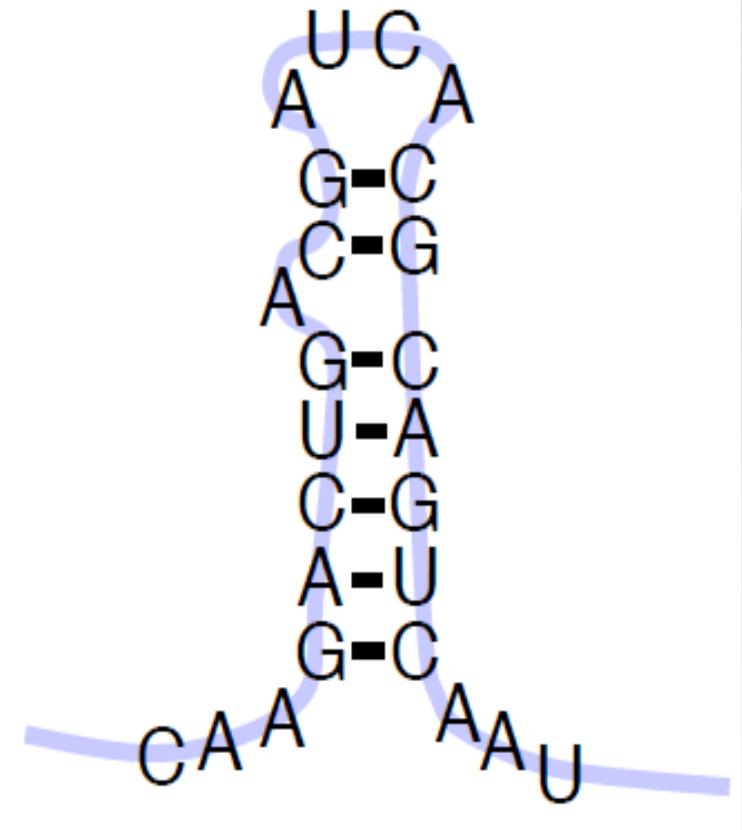
The RNA world

- An expending universe of RNA



→ Multiple roles of RNA in genes regulation

- RNA folds on itself by base pairing :
 - A with U : A-U, U-A
 - C with G : G-C, C-G
 - Sometimes G with U : U-G, G-U
- Folding = Secondary structure
- Structure related to function : ncRNA of the same family have a conserved structure
- Sequence less conserved

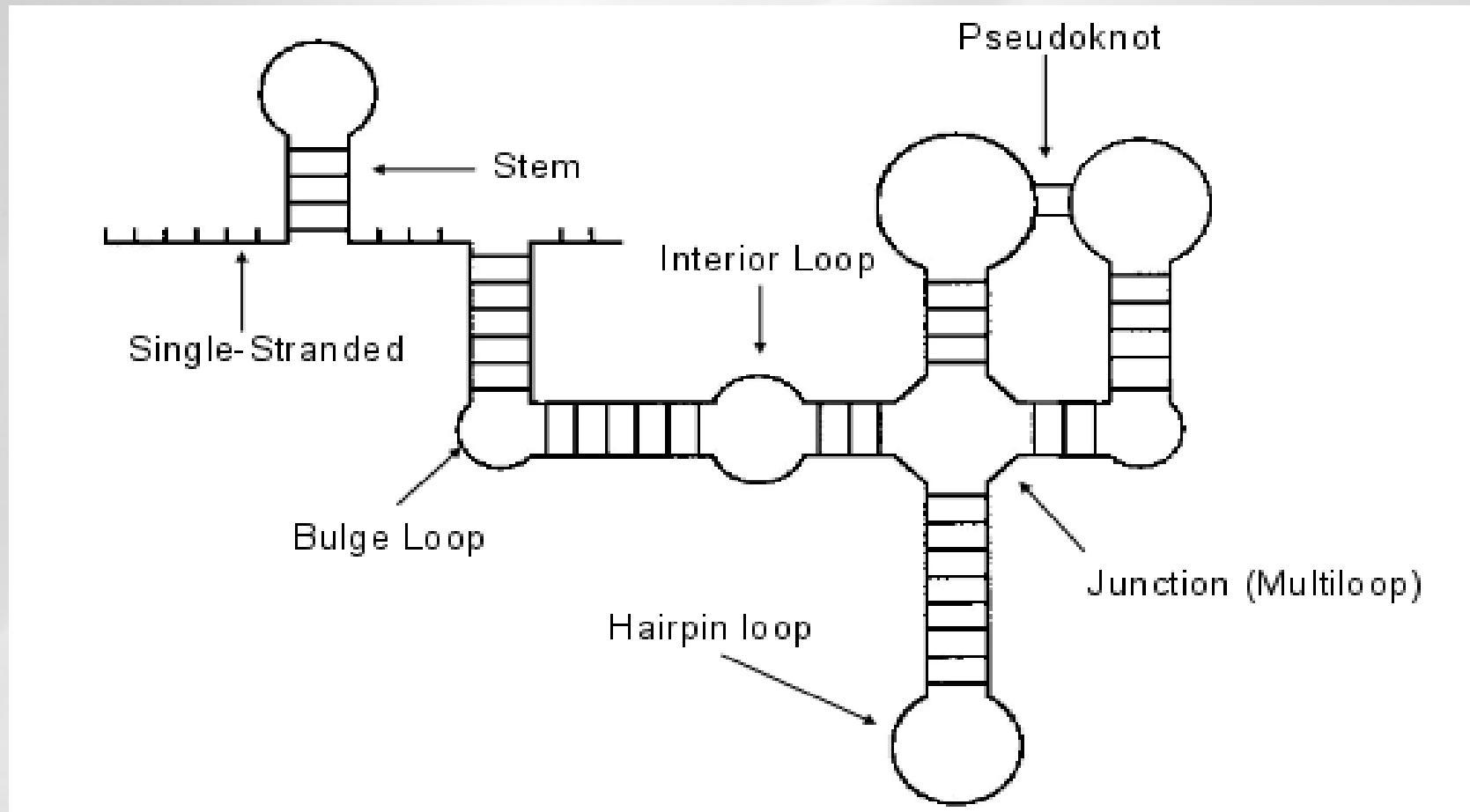


The non coding protein RNA world

- **Not predicted by gene prediction**
 - No specific signal (start, stop, splicing sites...)
 - Multiple location (intergenic, intronic, coding, antisens)
 - Variable size
 - No strong sequence conservation in general
- **A variety of existing approaches not always easy to integrate**
 - Known family: Homology prediction
 - New family: *De novo* prediction

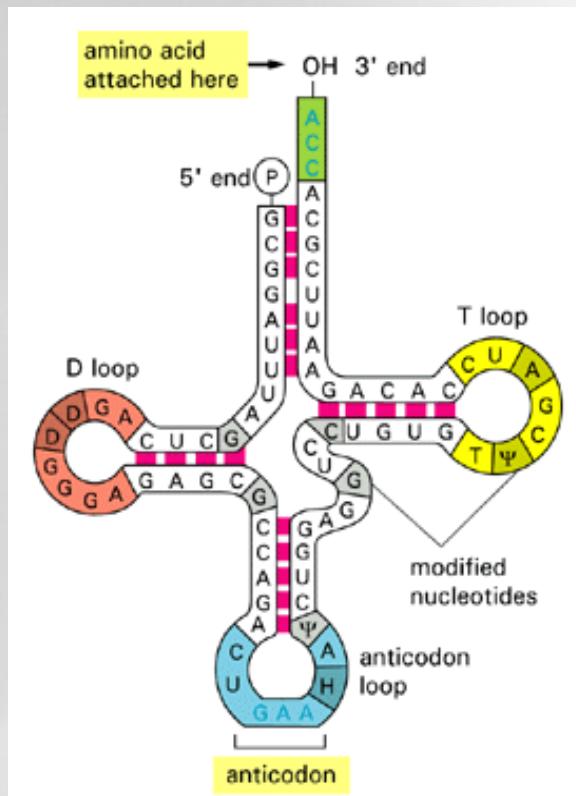
RNA background

Different elementary motifs



RNA background

Example: tRNA structure



The non coding protein RNA world

- **Large non coding protein RNA**
 - >300 nt
 - rRNA, tRNA, Xist, H19, ...
 - Genome structure & expression
- **Small non coding protein RNA**
 - >30 nt
 - snoRNA, snRNA...
 - mRNA maturation, translation
- **Micro non coding protein RNA**
 - 18-30 nt
 - miRNA, hc-siRNA, ta-siRNA, nat-siRNA, piRNA...
 - PTGS, TGS, Genome stability, defense...

Introduction to miRNA world and sRNAsseq

The non coding protein RNA world

- Large non coding protein RNA
 - >300 nt
 - rRNA, tRNA, Xist, H19, ...
 - Genome structure & expression
- Small non coding protein RNA
 - >30 nt
 - snoRNA, snRNA...
 - mRNA maturation, translation
- Micro non coding protein RNA
 - 18-30 nt
 - miRNA, hc-siRNA, ta-siRNA, nat-siRNA, piRNA...
 - PTGS, TGS, Genome stability, defense...

The miRNA world

• Discovery of *lin-4* in *C. elegans* in 1993



Cell, Vol. 75, 843–854, December 3, 1993, Copyright © 1993 by Cell Press

The *C. elegans* Heterochronic Gene *lin-4* Encodes Small RNAs with Antisense Complementarity to *lin-14*

Rosalind C. Lee,^{*†} Rhonda L. Feinb
and Victor Ambros^{*}
Harvard University
Department of Cellular and Developmental Biology
Cambridge, Massachusetts 02138

Summary

lin-4 is essential for the normal temporal pattern of diverse postembryonic development in *C. elegans*. *lin-4* acts by negatively regulating *lin-14* protein, creating a temporal gradient.

Cell, Vol. 75, 855–862, December 3, 1993, Copyright © 1993 by Cell Press

Posttranscriptional Regulation of the Heterochronic Gene *lin-14* by *lin-4* Mediates Temporal Pattern Formation in *C. elegans*

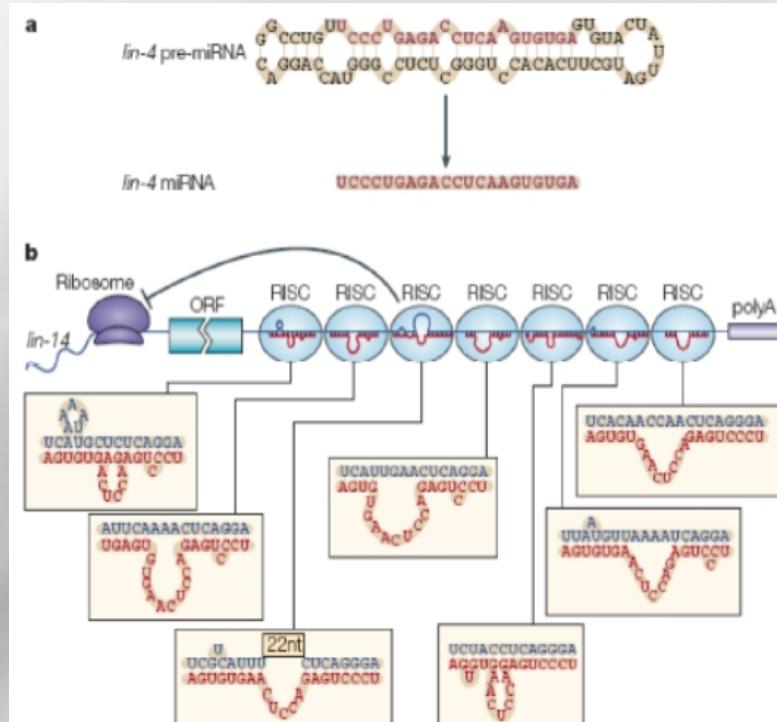
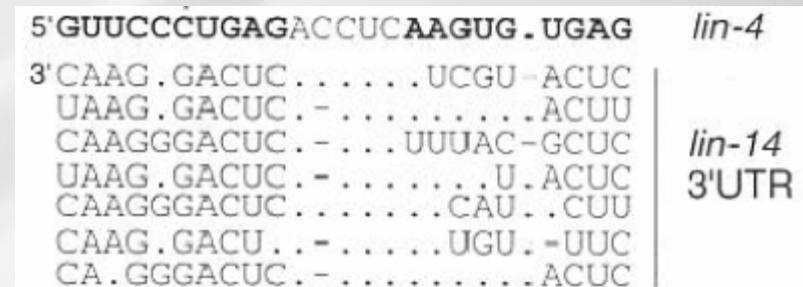
Bruce Wightman,^{*†} Ilho Ha,^{*} and Gary Ruvkun
Department of Molecular Biology
Massachusetts General Hospital
Boston, Massachusetts 02114

Summary

During *C. elegans* development, the temporal pattern of many cell lineages is specified by graded activity of the heterochronic gene *Lin-14*. Here we demonstrate

site phenotypes (Ambros and Horvitz, 1987). *lin-14(lf)* alleles cause larvae stage 2 (L2) patterns of cell lineage in a variety of tissues to be executed precociously during the L1 stage (Ambros and Horvitz, 1987). Two *lin-14(gf)* alleles cause the opposite transformation in temporal cell fate, reiterations of early cell fates at later stages. For instance, at the L2 stage, *lin-14(gf)* mutants repeat patterns of cell lineage appropriate for the L1 stage (Ambros and Horvitz, 1984).

lin-14 controls these stage-specific cell lineages by generating a temporal gradient of *Lin-14* nuclear protein (*Lin*



(He & Hannon, Nature reviews, 2004)

The miRNA world

- A key regulation function

Nature, 2011, January 20; 469(7330): 336–342, doi:10.1038/nature09783

Pervasive roles of microRNAs in cardiovascular biology

Eric M. Small¹ and Eric N. Olson¹

¹Department of Molecular Biology, University of Texas Southwestern Medical Center, Hines Boulevard, Dallas, Texas 75390-9148, USA

Development 138, 1081-1086 (2011) doi:10.1242/dev.056317
© 2011. Published by The Company of Biologists Ltd

Small RNAs Guide Hematopoiesis, Differentiation and Function

Francisco Navarro and Judy Lieberman

J Immunol 2010;184:5939-5947
doi:10.4049/jimmunol.0902567

<http://www.jimmunol.org/content/184> Byeong-Moo Kim^{1,2,*†}, Janghee Woo^{1,3,†}, Chryssa Kanelloupolou⁴ and Ramesh A. Shivdasani^{1,2,‡}

This information is current as of December 28, 2011

Developmental Cell 11, 441–450, October, 2006 © 2006 Elsevier Inc. DOI 10.1016/j.devcel.2006.09.001

379

The Diverse Functions of MicroRNAs in Animal Development and Disease

Wigard P. Kloosterman¹ and Ronald H.A. Plasterk^{1,2,*}

¹ Hubrecht Laboratory
Centre for Biomedical Genetics

Leading Edge
Review

Origin, Biogenesis, and Activity of Plant MicroRNAs

Olivier Voynette^{1,*}

Olivier Voynnet*
Institut de Biologie Moléculaire des Plantes, CNRS UPR2357–Université de Strasbourg, 67084 Strasbourg, France
*Correspondence: olivier.voynnet@ibmp-ulp.u-strasbg.fr
DOI 10.1016/j.cell.2009.01.046

MicroRNAs (miRNAs) are key posttranscriptional regulators of eukaryotic gene expression. They use highly conserved as well as more recently evolved, species-specific mechanisms to regulate an array of biological processes. This Review discusses current advances in our understanding of miRNA origin, biogenesis, and mode of action of plant miRNAs and draws comparisons with their animal counterparts.

Since then, several RNA-cloning strategies to vertebrates and invertebrates

The discovery of hundreds of plant micro RNAs (miRNAs) has triggered much speculation about their potential roles in plant development. The search for plant genes involved in miRNA processing has revealed common factors such as DICER, and new molecules, including HEN1. Progress is also being made toward identifying miRNA target genes and understanding the mechanisms of miRNA-mediated gene regulation in plants. This work has lead to a reexamination of n characterized mutations that are now International Jour-

Addresses

This review comes from a themed issue

Pattern formation and developmental
Edited by Anne Echard and Olivier R.

International Journal of Alzheimer's Disease
Volume 2011 (2011), Article ID 894938, 6 pages
<https://doi.org/10.4061/2011/894938>

Review Article

MicroRNAs and Alzheimer's Disease Mouse Models: Current Insights and Future Research Avenues

Charlotte Delay^{1,2} and Sébastien S. Hébert^{1,2}

The miRNA world

Regulatory functions

• Animals

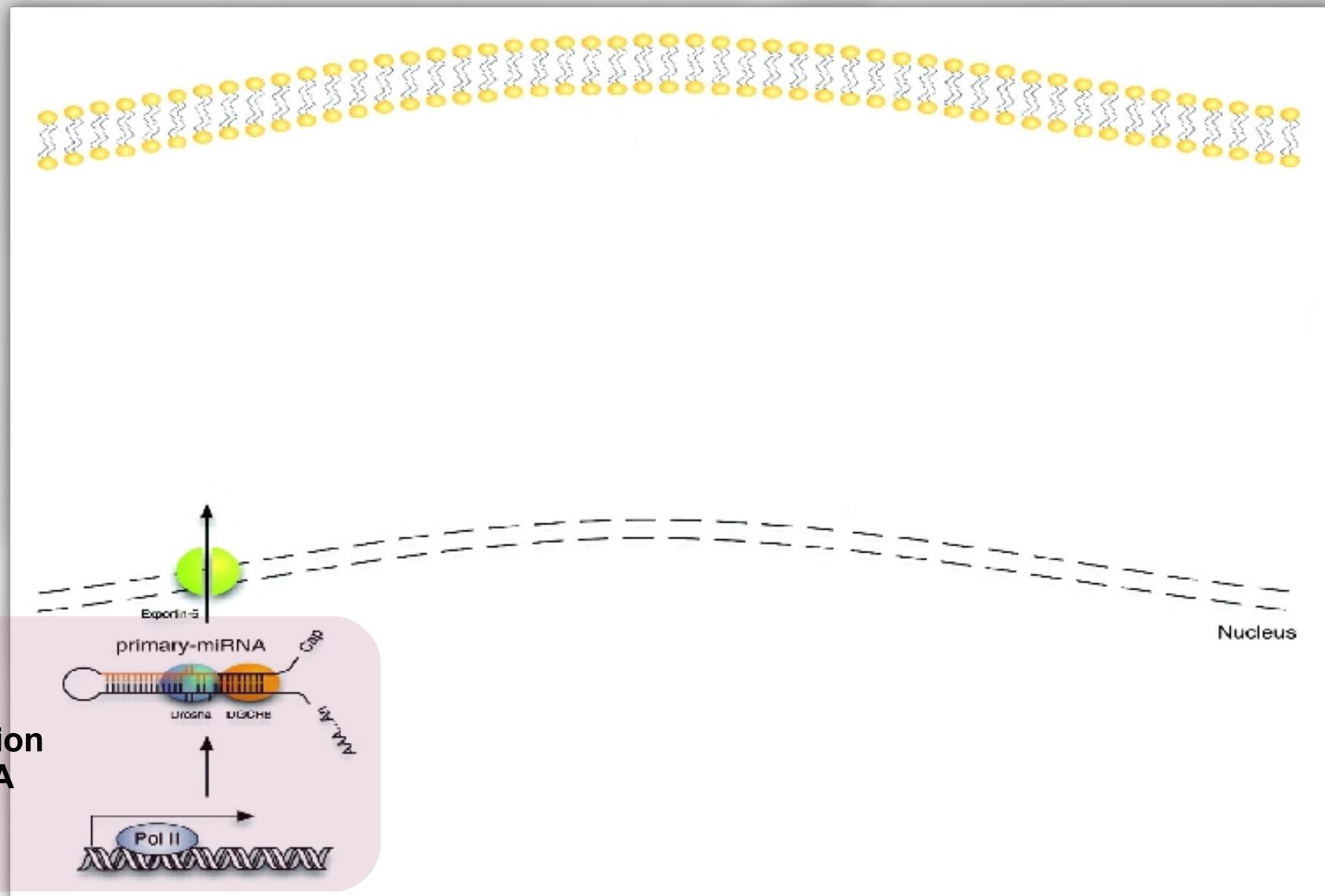
- Developmental timing (*C. elegans*): **lin-4, let-7**
- Neuronal left/right asymmetry (*C. elegans*): **Lys-6, mir-273**
- Programmed cell death/fat metabolism (*D. melanogaster*): **mir-14**
- Notch signaling (*D. melanogaster*): **mir-7**
- Brain morphogenesis (Zebrafish): **mir-430**
- Myogeneses and cardiogenesis: **mir-1, miR-181, miR-133**
- Insulin secretion: **miR-375**
- ...

1600 precursors in Human !!! (ref: miRBase, August 2012)

• Plants

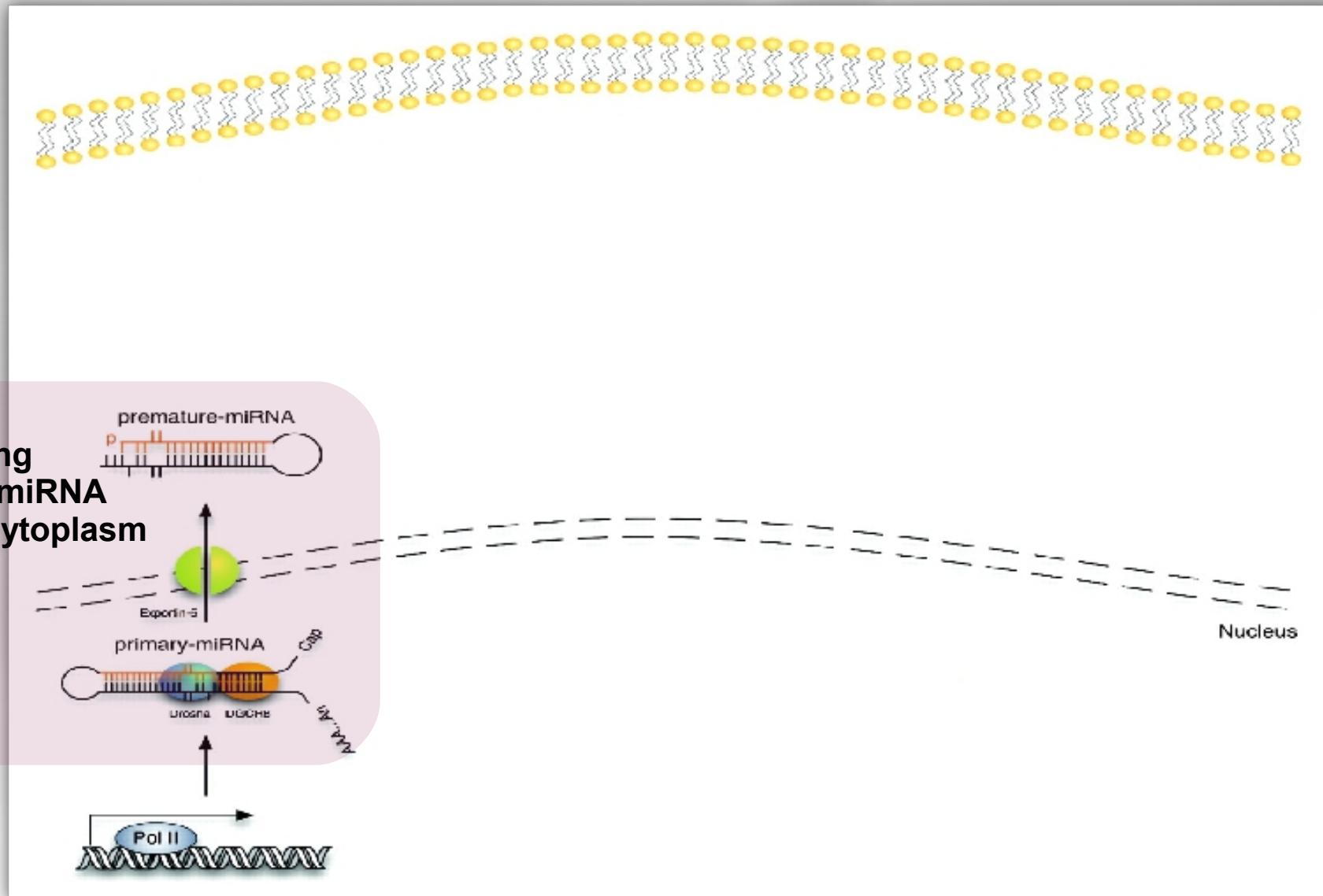
- Floral timing and leaf development: **miR-156**
- Organ polarity, vascular and meristem development: **mir-165, miR-166**
- Expression of auxin response genes: **miR-160**
- ...

The miRNA biogenesis



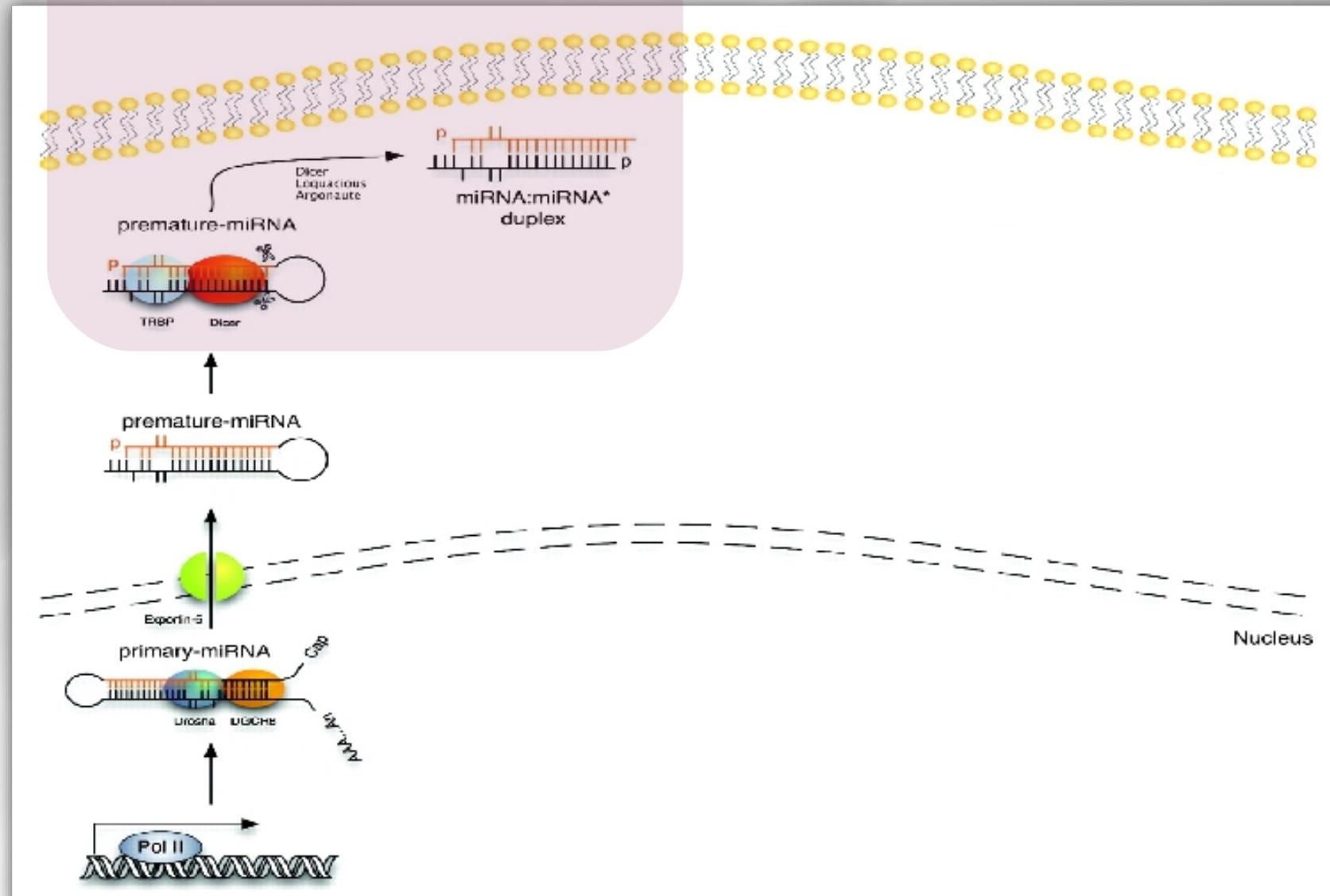
Pol II transcription
Into a pri-miRNA

The miRNA biogenesis

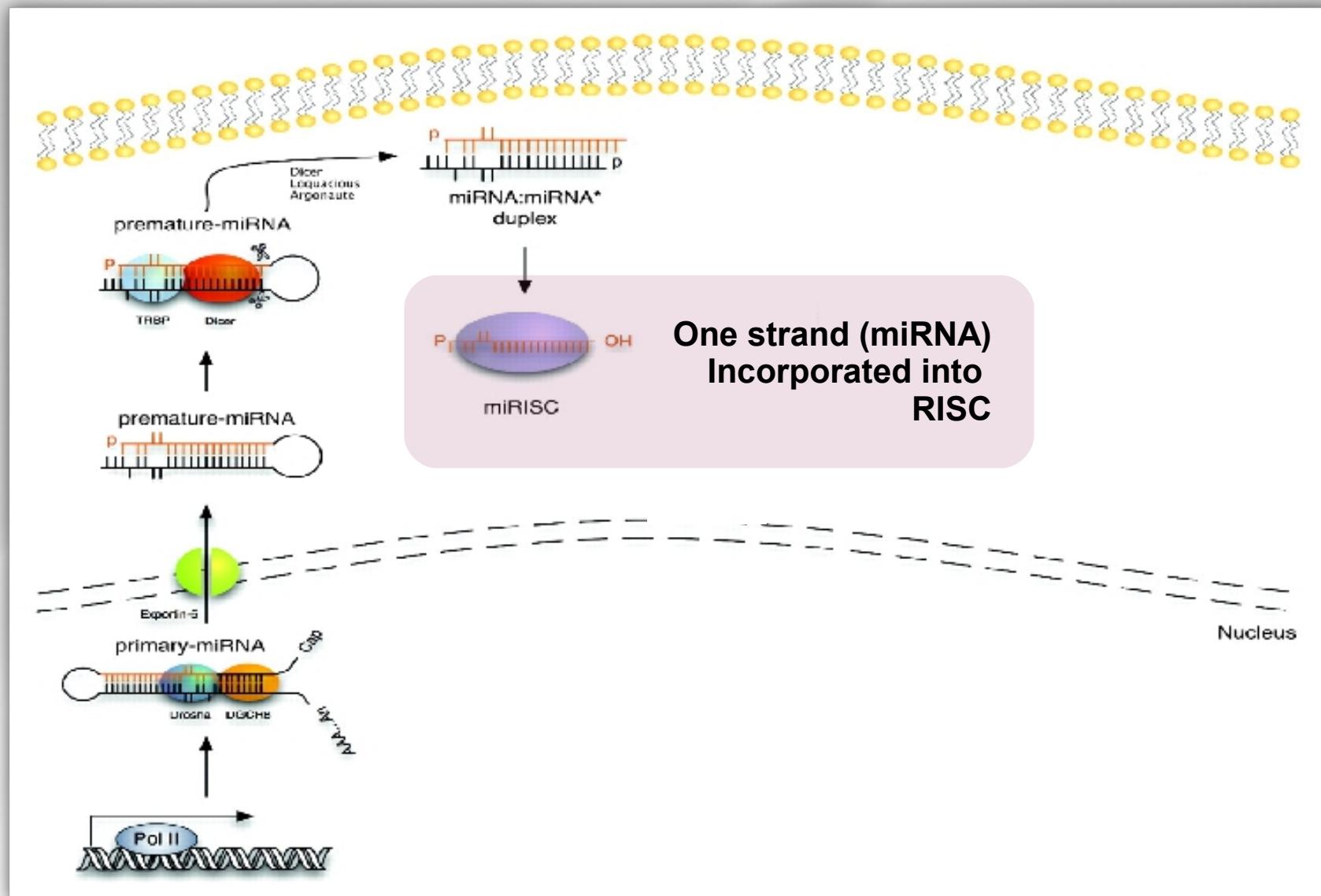


The miRNA biogenesis

Dicer processing Into a duplex miRNA Structure

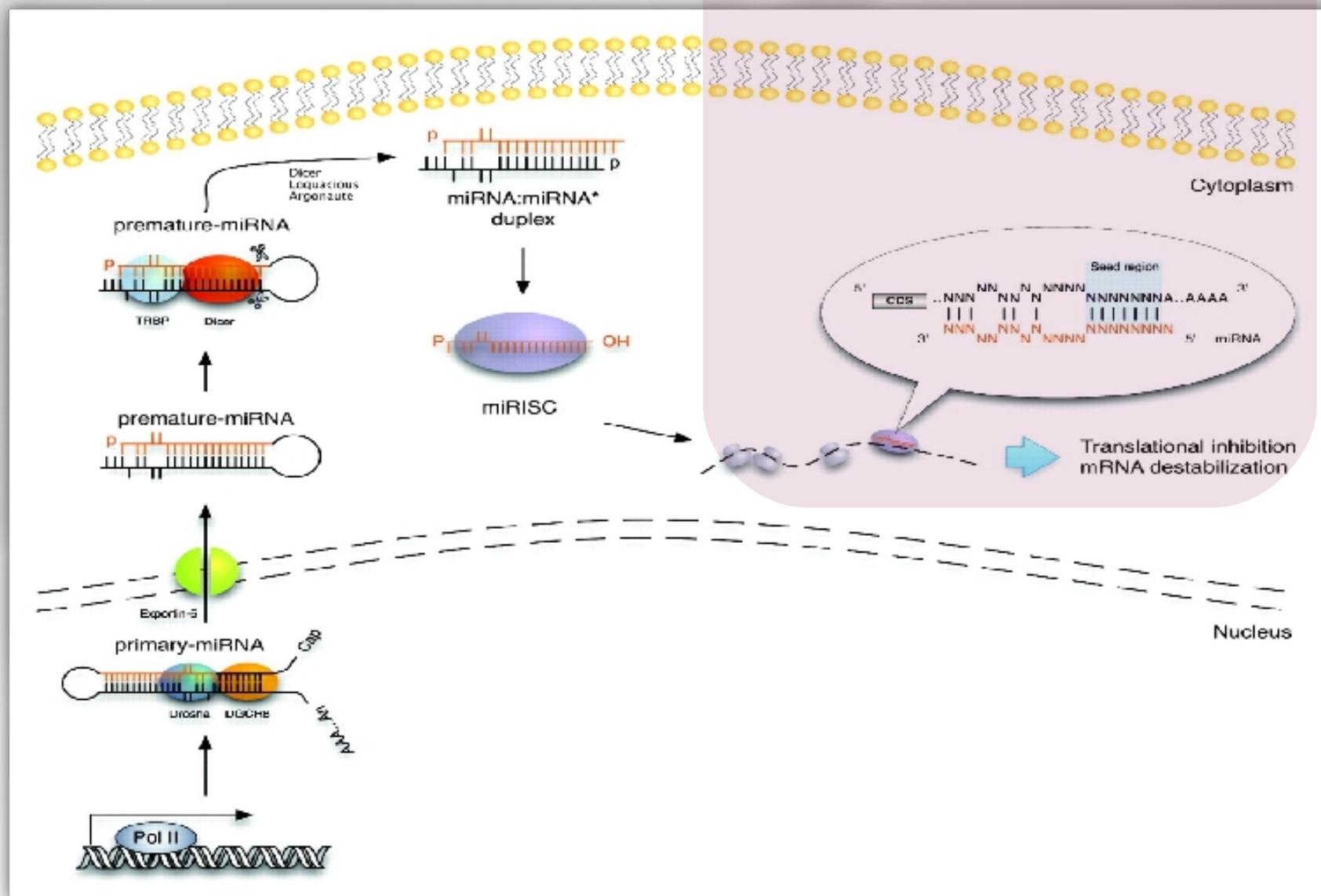


The miRNA biogenesis



The miRNA biogenesis

target mRNA
translationally
repressed



The Mirtron Pathway Generates microRNA-Class Regulatory RNAs in *Drosophila*

Katsutomo Okamura,¹ Joshua W. Hagen,¹ Hong Duan,¹ David M. Tyler,¹ and Eric C. Lai^{1,*}
¹Memorial Sloan-Kettering Cancer Center, Department of Developmental Biology, 1275 York Ave, Box 252, New York, NY 10021, USA
*Correspondence: laie@mskcc.org

Molecular Cell

Volume 28, Issue 2, 26 October 2007, Pages 328–336

Resource

Mammalian Mirtron Genes

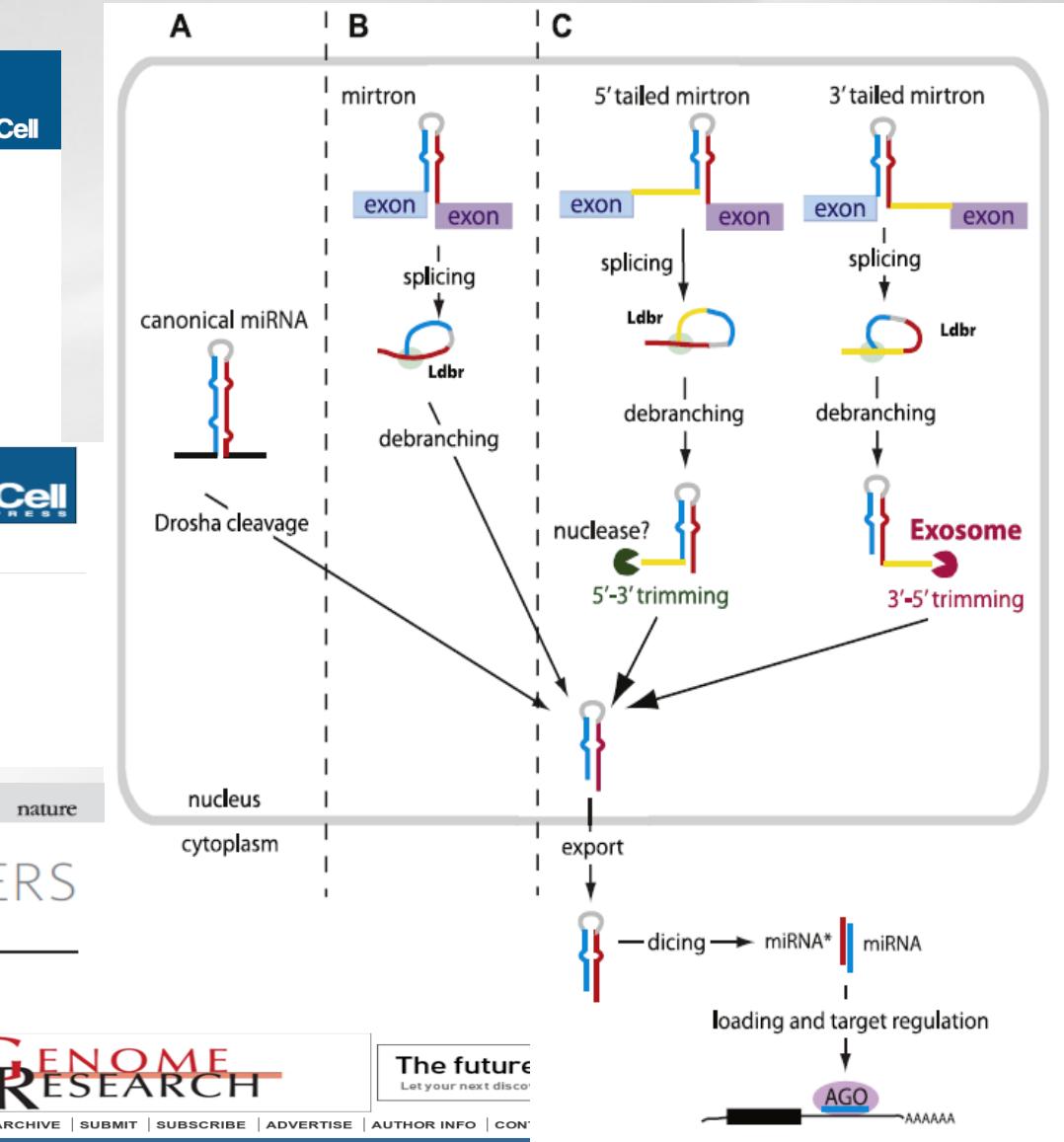
Eugene Berezikov¹,    Wei-Jen Chung², Jason Willis², Edwin Cuppen¹, Eric C. Lai²,   

¹ Hubrecht Institute, Uppsalalaan 8, 3584 CT Utrecht, The Netherlands

² Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, NY 10021, USA

<http://dx.doi.org/10.1016/j.molcel.2007.09.028>, How to Cite or Link Using DOI

Vol 448 | 5 July 2007 | doi:10.1038/nature05983



Intronic microRNA precursors that bypass Drosha processing

J. Graham Ruby^{1,2*}, Calvin H. Jan^{1,2*} & David P. Bartel^{1,2}



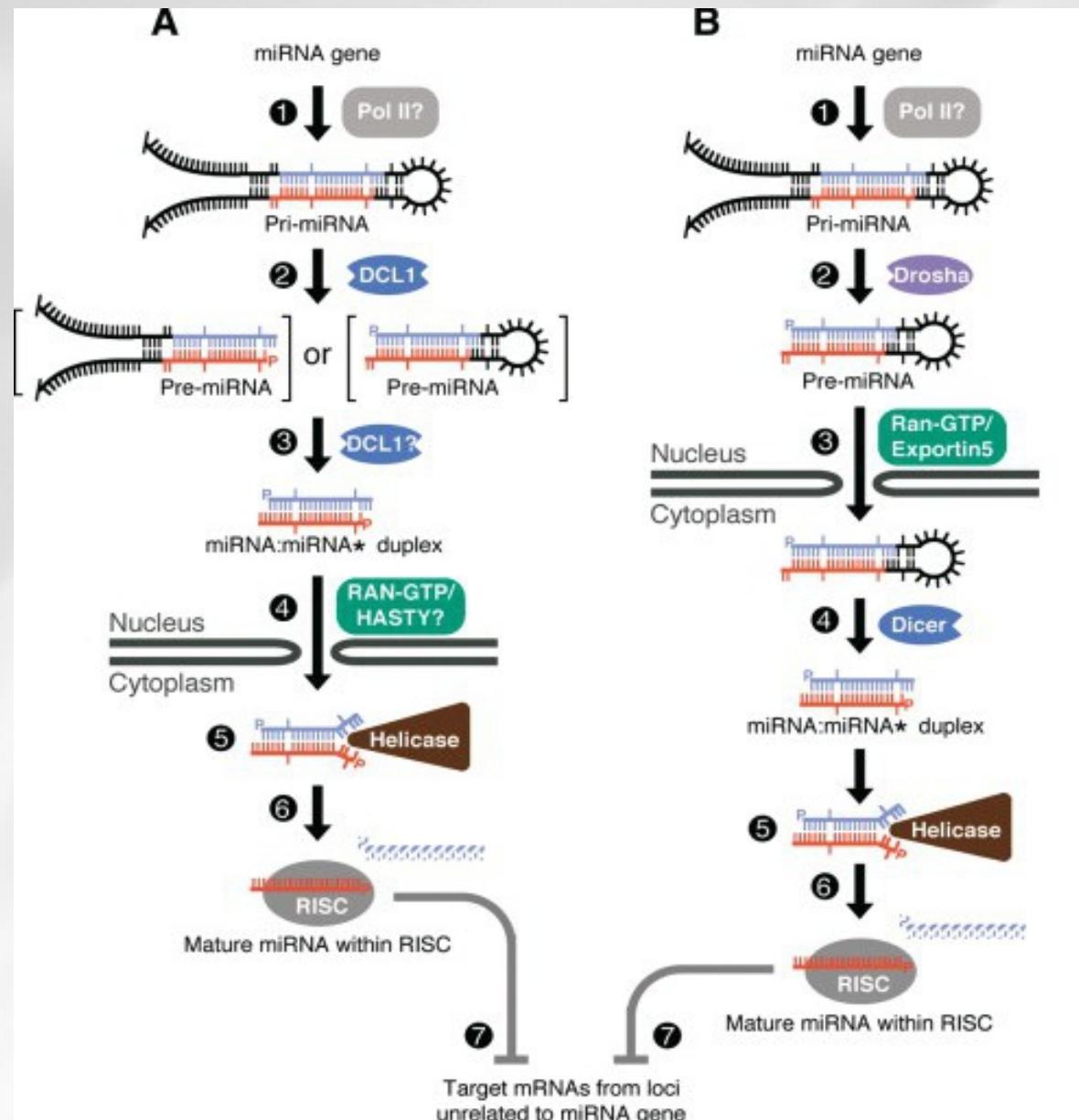
The future
Let your next discov

Institution: INRA Institut National de la Recherche Agronomique Sign In via U

Discovery of hundreds of mirtrons in mouse and human small RNA data

Erik Ladewig¹, Katsutomo Okamura^{1,2}, Alex S. Flynt¹, Jakub O. Westholm¹ and Eric C. Lai^{1,3}

The miRNA biogenesis



The miRNA location

a Non-coding TU with intronic miRNA

DLEU2



b Non-coding TU with exonic miRNA

B1C



c Coding TU with intronic miRNA

MCM7



d Coding TU with exonic miRNA

CACNG8



→ Cluster organisation

miRNA conservation

Conservation of the sequence and temporal expression of *let-7* heterochronic regulatory RNA

Amy E. Pasquinelli^{*†}, Brenda J. Reinhart^{*†}, Frank Slack[‡],
 Mark Q. Martindale[§], Mitzi I. Kurodall, Betsy Maller[‡], David C. Hayward[¶],
 Eldon E. Ball[¶], Bernard Degnan[#], Peter Müller[★], Jürg Spring[★],
 Ashok Srinivasan^{**}, Mark Fishman^{**}, John Finnerty^{††}, Joseph Corbo^{‡‡},
 Michael Levine^{‡‡}, Patrick Leahy^{§§}, Eric Davidson^{§§} & Gary Ruvkun^{*}

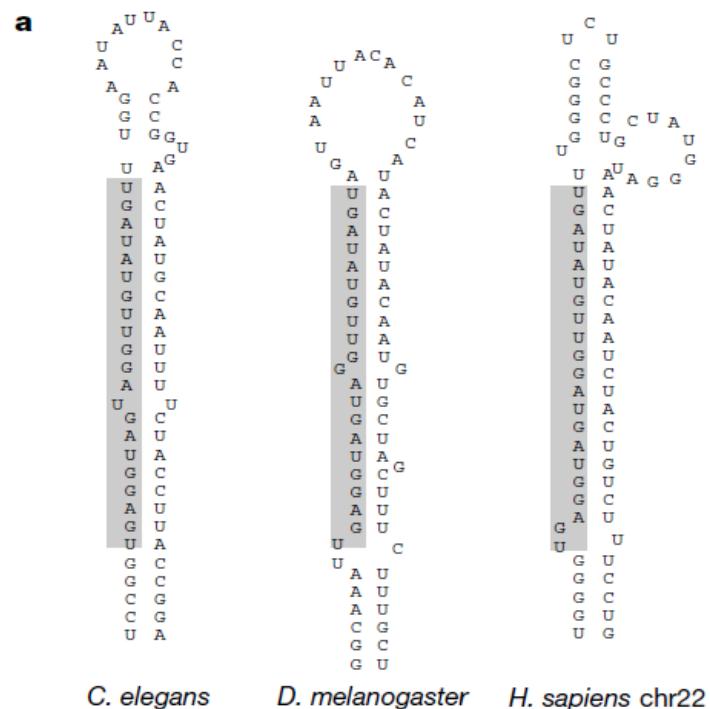
^{*} Department of Molecular Biology, Massachusetts General Hospital, and
 Department of Genetics, Harvard Medical School, Boston, Massachusetts 02114,
 USA

[†] Department of Molecular, Cellular and Developmental Biology, Yale University,
 New Haven, Connecticut 06520, USA

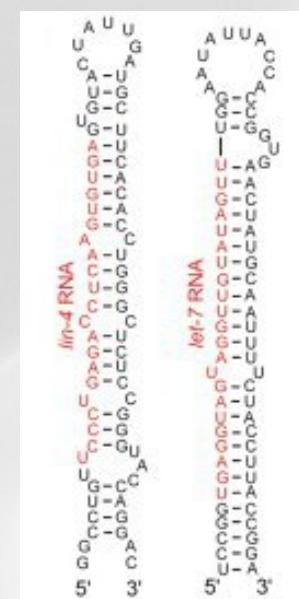
[‡] Kewalo Marine Lab, Pacific Biomedical Research Center, University of Hawaii,
 Honolulu, Hawaii 96813, USA

[¶] Howard Hughes Medical Institute, Baylor College of Medicine, Houston,
 Texas 77030, USA

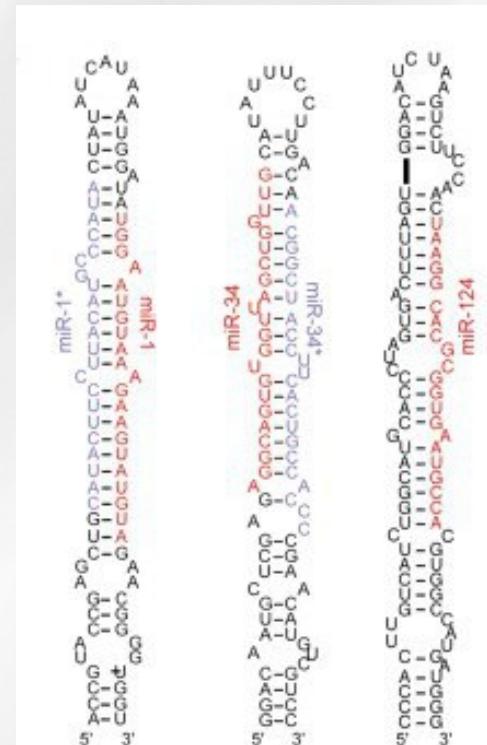
© 2000 Nature Publishing Group



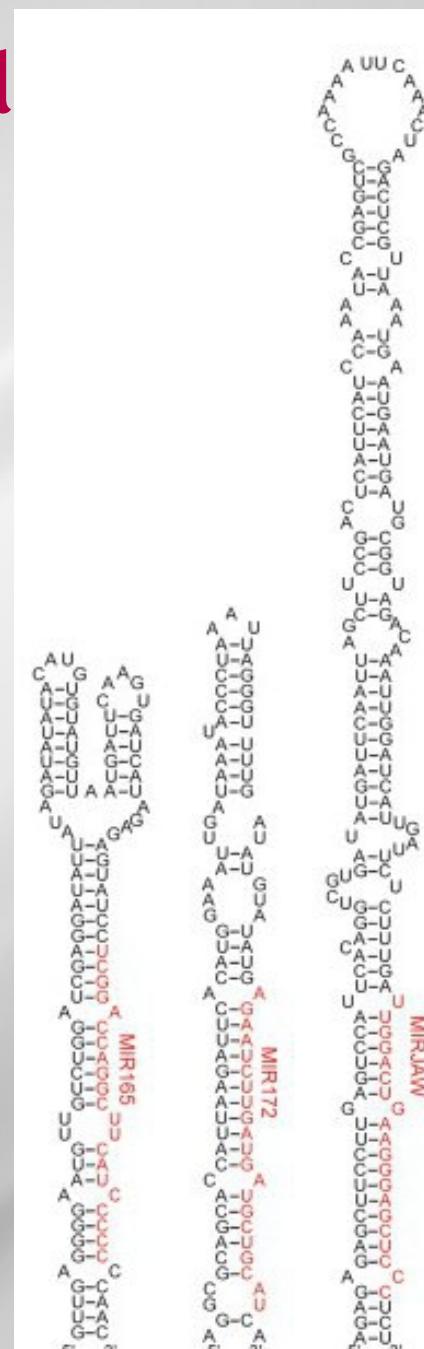
A. E. Pasquinelli et al., Nature 408, 86-9 (2000)



Initial miRNA



Animal miRNA



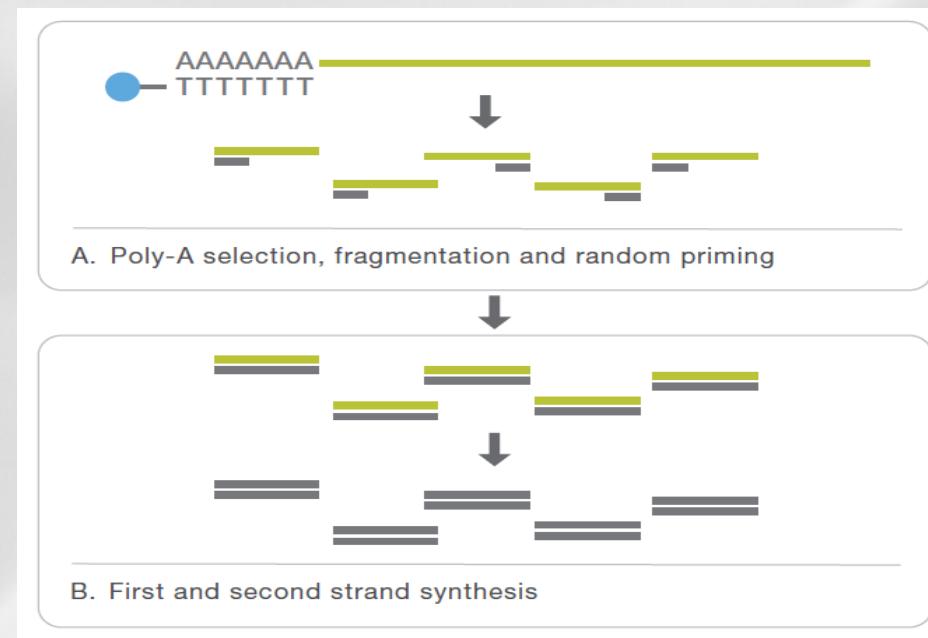
Plant miRNA

miRNA: pl

s animal

How can we study miRNA ?

- RNAseq not suited for miRNA (protocol and size)



- small RNAseq: ability of high throughput sequencing to
 - Interrogate known and new small RNAs
 - Quantify them
 - Profile them on a large number of samples
 - Cost-effective

small RNAseq platforms comparisons

OPEN  ACCESS Freely available online

Deep-Sequencing Protocols Influence the Results Obtained in Small-RNA Sequencing

Joern Toedling^{1,2,3,4,5*}, Nicolas Servant^{1,2,3*}, Constance Ciaudo^{1,4,5,6*}, Laurent Farinelli⁷, Olivier Voinnet^{6,8}, Edith Heard^{1,4,5†}, Emmanuel Barillot^{1,2,3†}

1 Institut Curie, Paris, France, **2** INSERM U900, Paris, France, **3** Mines ParisTech, Fontainebleau, France, **4** CNRS UMR3215, Paris, France, **5** INSERM U934, Paris, France,

6 Department of Biology, Swiss Federal Institute of Technology Zürich, Zürich, Switzerland, **7** Fasteris, Plan-les-Ouates, Switzerland, **8** Institut de Biologie Moléculaire des Plantes, CNRS UPR2357 – Université Louis Pasteur, Strasbourg, France

Table 1. Description of the libraries investigated.

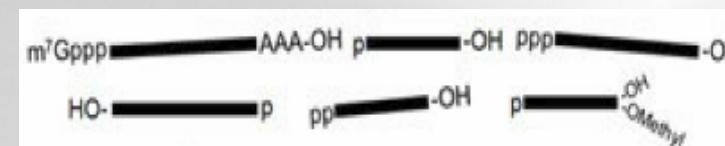
SampleID	CellType	Technology	Year	Barcode/index	Comment	# reads
ES_XY_454	E14 XY	454	2008	barcode	Ciaudo et al. (2009)	95203
ES_XY_Solexa_Illu	E14 XY	Solexa	2010	none	GAIix/Illumina 3' adapter	28014973
ES_XY_Solexa_i_IDT	E14 XY	Solexa	2010	Index	HiSeq2000/IDT 3' adapter	8375905
ES_XY_Solexa_IDT	E14 XY	Solexa	2010	none	GAIix/IDT 3' adapter	31316082
ES_XY_SOLiD_v3	E14 XY	SOLiD	2010	none	v3+/SREK kit	32685742
ES_XY_SOLiD_v4	E14 XY	SOLiD	2010	barcode	v4/STaR-Seq kit	2685423
ES_XX_454	PGK XX	454	2008	barcode	Ciaudo et al. (2009)	57497
ES_XX_Solexa_i_IDT	PGK XX	Solexa	2009	Index	HiSeq2000/IDT 3' adapter	10262556
ES_XX_SOLiD_v3	PGK XX	SOLiD	2010	none	v3+/SREK kit	32974547
ES_XX_SOLiD_v4	PGK XX	SOLiD	2010	barcode	v4/STaR-Seq kit	2714593

The ten samples differ in size, in the employed sequencing technology, in the version of the machine that they were generated with and whether a barcode or index had been used for parallel sequencing with other libraries.

doi:10.1371/journal.pone.0032724.t001

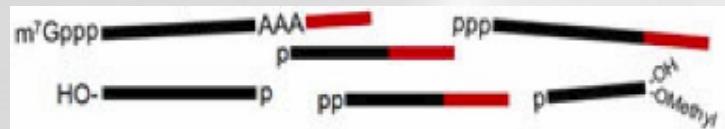
small RNA-Seq library preparation

- Monophosphate presence in 5' extremity and OH presence in 3' extremity



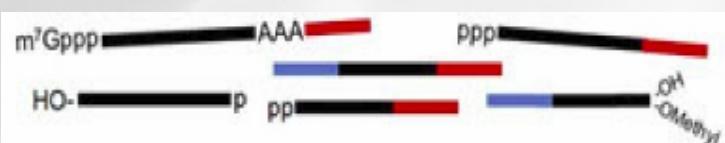
Total RNA: contain all kinds of RNA species including miRNA, mRNA, tRNA, rRNA...

↓ **Ligate with 3' adapter**



RNA with modified 3'-end will not ligate with 3' adapters. Only RNA with OH in 3'-end will ligate.

↓ **Ligate with 5' adapter**



Only RNA with monophosphate in 5'-end will ligate with 5' adapters.

↓ **RT-PCR and Size Selection**



MicroRNA sequencing library

CDNA containing both adapter sequences will be amplified. MicroRNA will be enriched from PCR and gel size selection.

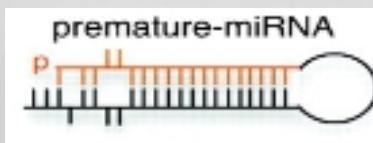
What are we looking for ?

- **List of known miRNA**
- **List of new miRNA**
- **miRNA target(s)**
- **miRNA quantification**
- **Differential expression**

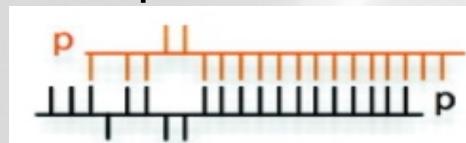
small RNAseq data analysis

What should we retain for data analysis ?

- Pre-miRNA information:



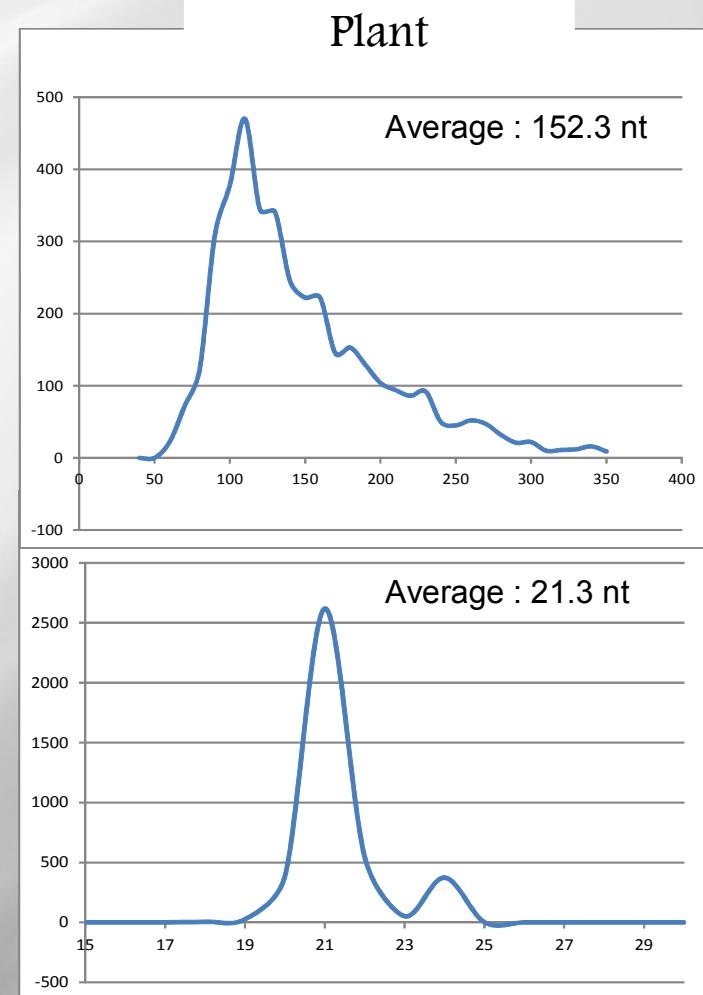
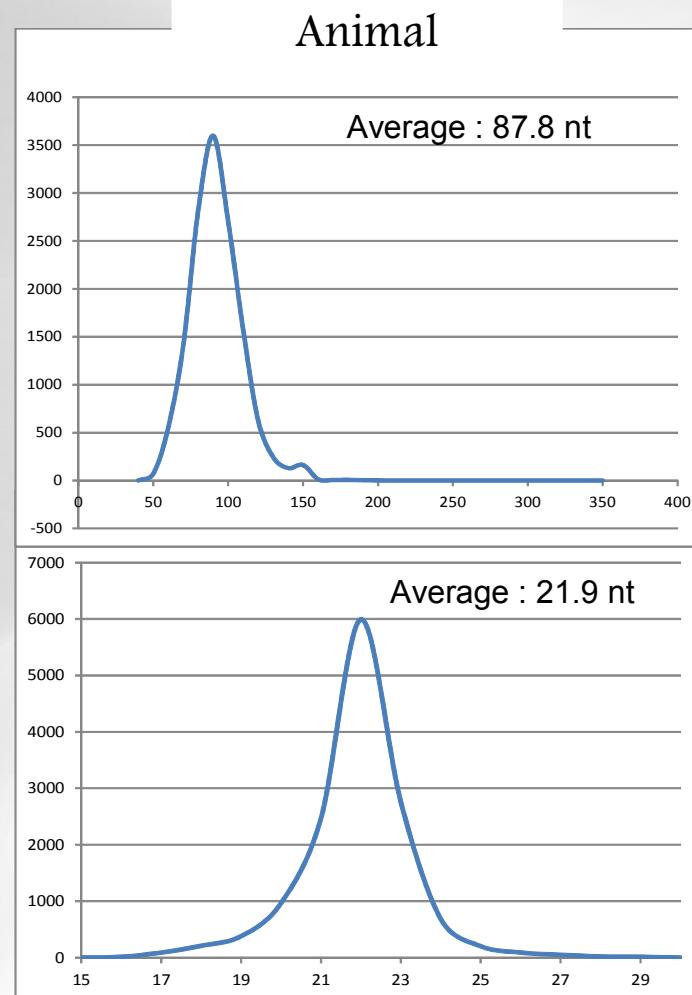
- Hairpin structure of the pre-miRNA
- Pre-miRNA localisation (coding/non coding TU intronic/exonic)
- Presence of cluster
- Size of the pre-miRNA
- miRNA-5p and miRNA-3p information:



- Existence of both miRNA-5p and miRNA-3p
- Sequence conservation
- Overhang (around 2 nt) related to drosha and Dicer cuts
- Size of miRNA-5p and miRNA-3p
- Overexpression of one of the miRNA-5p and miRNA-3p
- Existence of other products in sRNAseq data

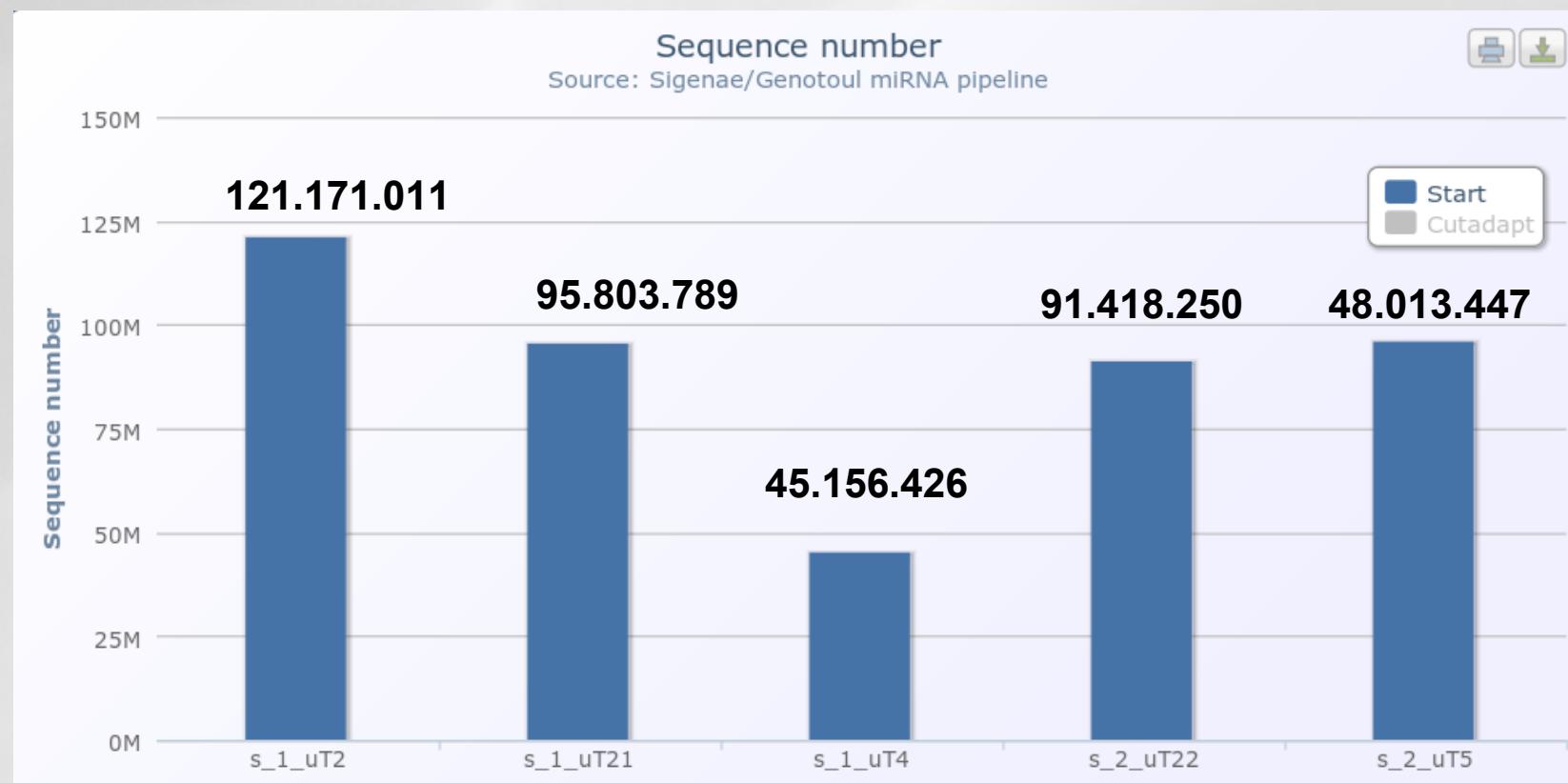
What should we retain for data analysis ?

miRbase data on pre-miRNA / mature



Description of the dataset

- 5 experiments (5 lanes, no multiplexing)
 - Different tissues, different stages
- No reference genome
 - Only scaffolds



Fastq format

@D61655M1_171:2:1:1192:1017#0/1
NN
+D61655M1_171:2:1:1192:1017#0/1
BB
@D61655M1_171:2:1:1202:1038#0/1
NN
+D61655M1_171:2:1:1202:1038#0/1
BB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTGCTTGAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y [YXXccccccccc\cccc_aacccYUUUVV0Q
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTGCTTGAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V\^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGGCATCTCGTATGCCGT
+D61655M1_171:2:1:13819:1998#0/1
ggggggggggfgfgggf^ggggfggggegggdgggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTGTCAGACTTTGTTGGCATCTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB

Fastq format

Line 1 starts with @

Information	Meaning
D61655M1_171	The unique instrument name
2	Flowcell lane45.156.426
1	Tile number within the flow cell lane
1192	'x'-coordinate of the cluster within the tile
1017	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 (<i>paired-end or mate-pair reads only</i>)

Fastq format

Line 1 starts with @

Information	Meaning
D61655M1_171	The unique instrument name
2	Flowcell lane45.156.426
1	Tile number within the flow cell lane
1192	'x'-coordinate of the cluster within the tile
1017	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 (<i>paired-end or mate-pair reads only</i>)

Line 2 Raw sequence of 36 nt (36 cycles in sequencing)

Fastq format

Line 1 starts with @

Information	Meaning
D61655M1_171	The unique instrument name
2	Flowcell lane45.156.426
1	Tile number within the flow cell lane
1192	'x'-coordinate of the cluster within the tile
1017	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 (<i>paired-end or mate-pair reads only</i>)

Line 2 Raw sequence of 36 nt (36 cycles in sequencing)

Line 3 starts with a '+' character and is optionally followed by the same sequence identifier (and any description) again.

Fastq format

Line 1 starts with @

Information	Meaning
D61655M1_171	The unique instrument name
2	Flowcell lane45.156.426
1	Tile number within the flow cell lane
1192	'x'-coordinate of the cluster within the tile
1017	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 (<i>paired-end or mate-pair reads only</i>)

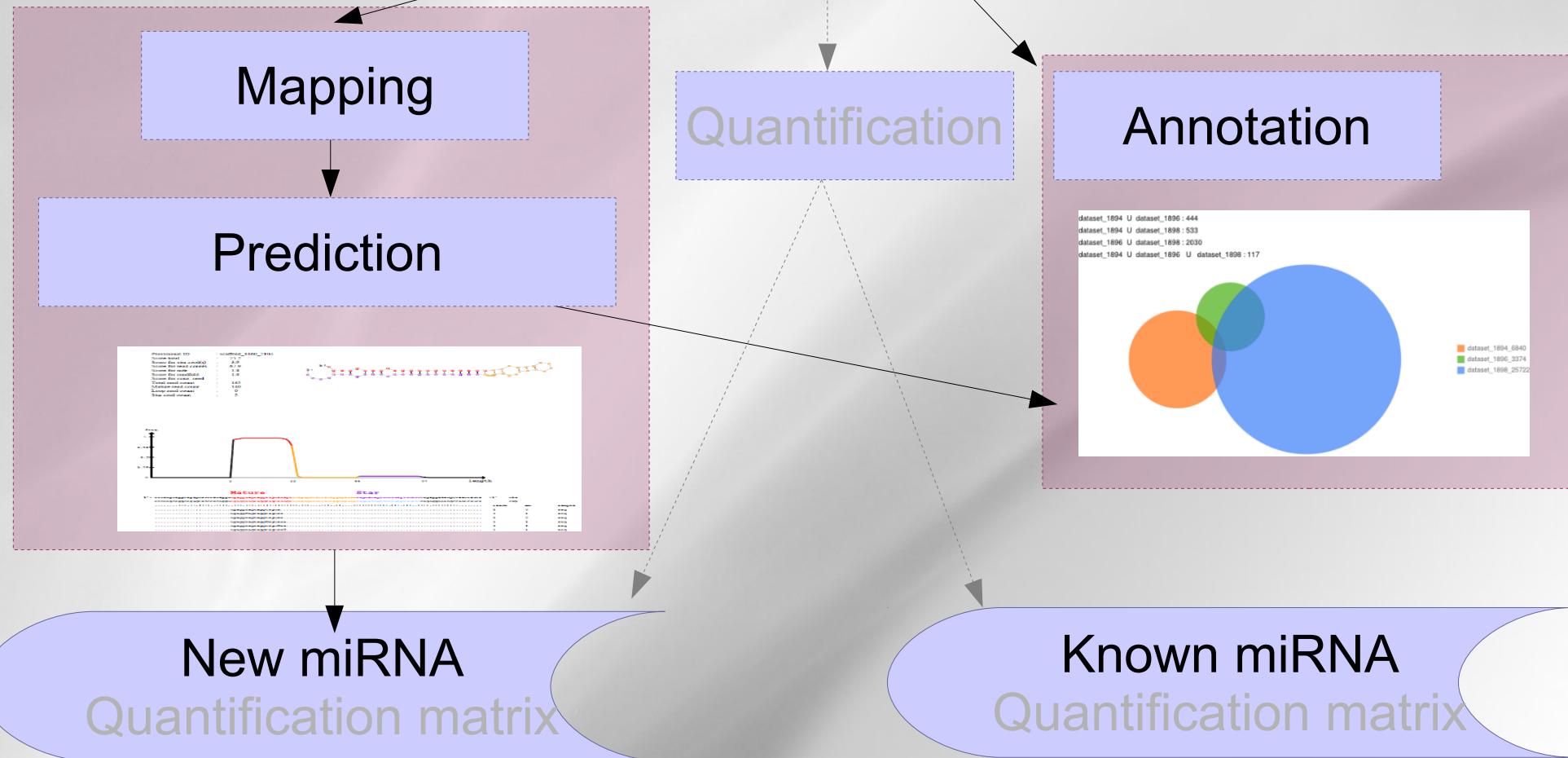
Line 2 Raw sequence of 36 nt (36 cycles in sequencing)

Line 3 starts with a '+' character and is optionally followed by the same sequence identifier (and any description) again.

Line 4 Line 4 encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence.

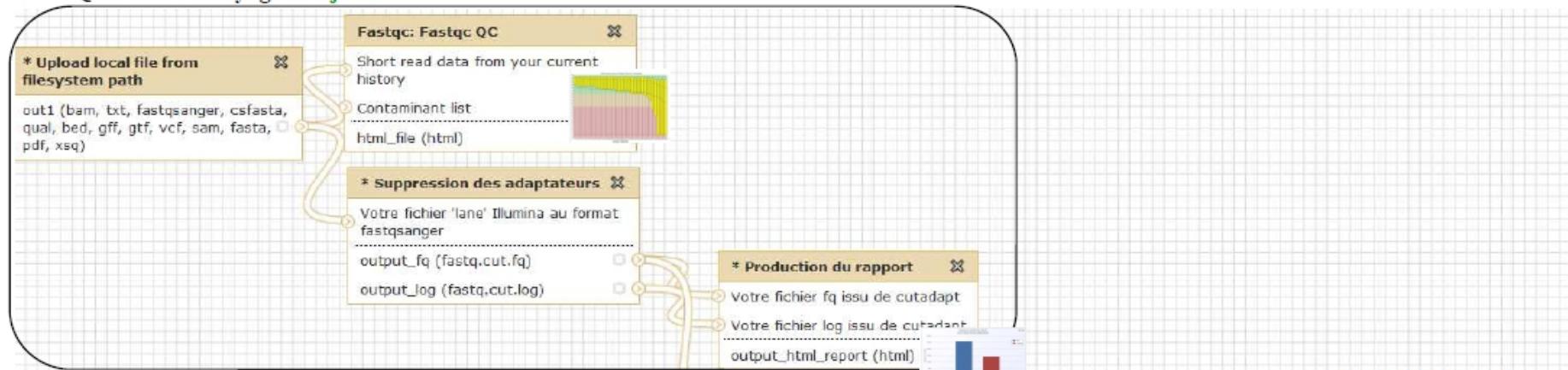
small RNAseq pipeline

with reference



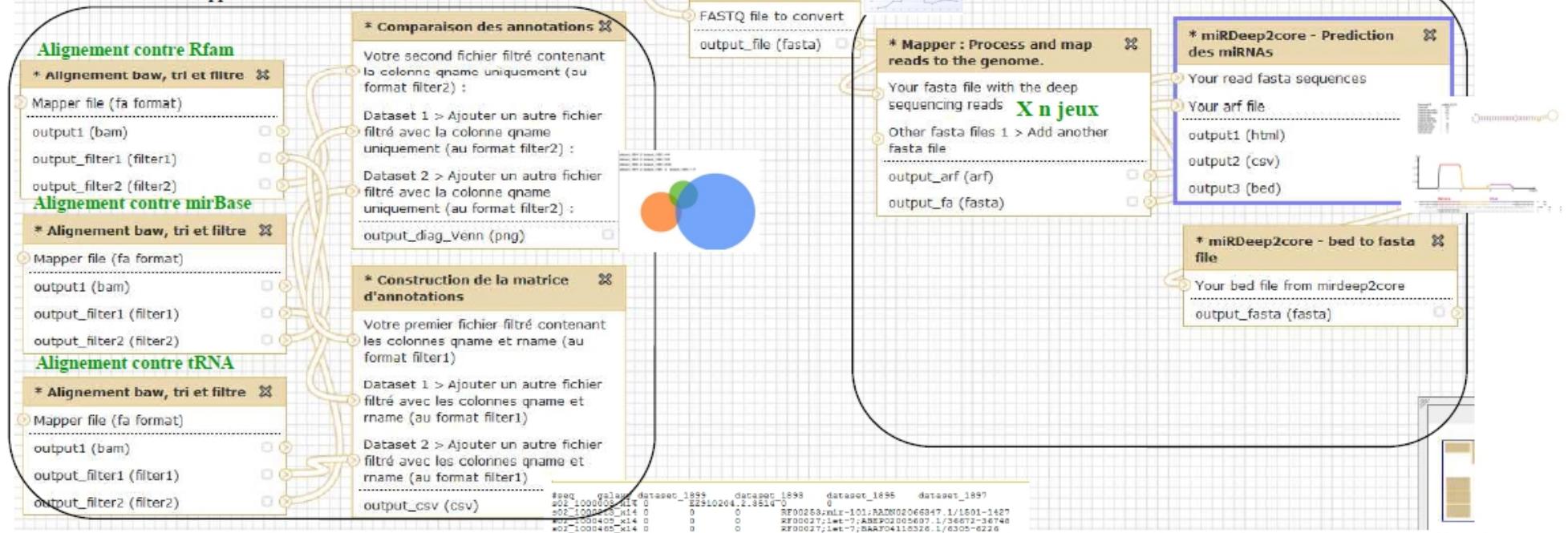
small RNAseq pipeline

WF1 Qualité et nettoyage X n jeux



WF3 Annotations fonctionnelles

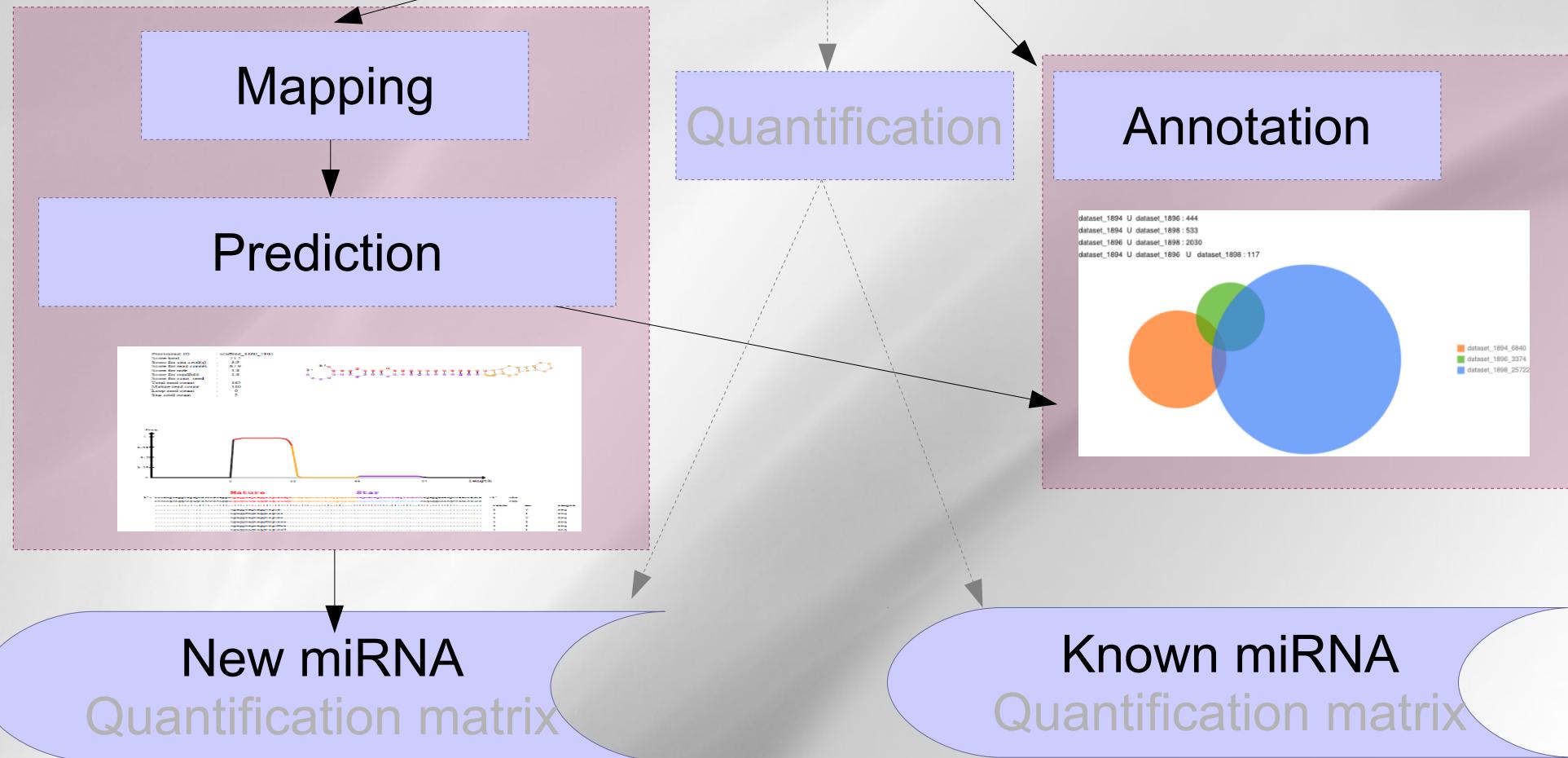
Avec le fasta issu de « mapper » et/ou le fasta issu de « core »



seq	qual	dataset_1899	dataset_1899	dataset_1899	dataset_1899	dataset_1899
s02-10000000x14	0	0	0	0	R000259:mnr-101:RANB02068347.1/1501-1427	
s02-10000001x14	0	0	0	0	R000257:lat-7:RANP02005607.1/36672-36740	
s02-10000002x14	0	0	0	0	R000257:lat-7:RANP04115328.1/6305-6226	
s02-10000003x14	0	0	0	0	R000257:lat-7:RANP04115328.1/4075-5290	
s02-10000004x14	0	0	0	0	R000257:lat-7:RANP04115328.1/1844-18249	
s02-10000005x14	0	0	0	0	FJ698877.1.3204	0
s02-10000006x14	0	0	0	0	R000257:lat-7:RANP04115328.1/18701-18715	
s02-10000007x14	0	0	0	0	R000257:lat-7:RANP04115328.1/18725-18745	
s02-10000008x14	0	0	0	0	R000259:mnr-180:RANB01010010.1/106564-10647	
s02-10000009x14	0	0	0	0	EU875689.103747	118671 0
s02-10000010x14	0	0	0	0	R000257:lat-7:RANP04096587.1/14846-14268	
s02-10000011x14	0	0	0	0	R000257:lat-7:RANP04096587.1/21463-21220	

small RNAseq pipeline

with reference



- **FastQC** (<http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/>)

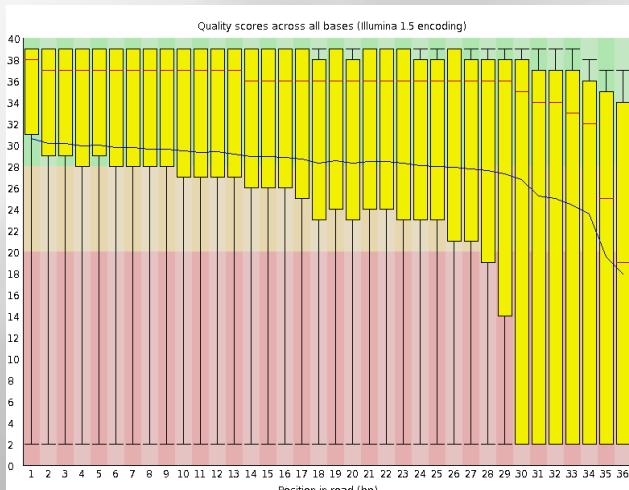
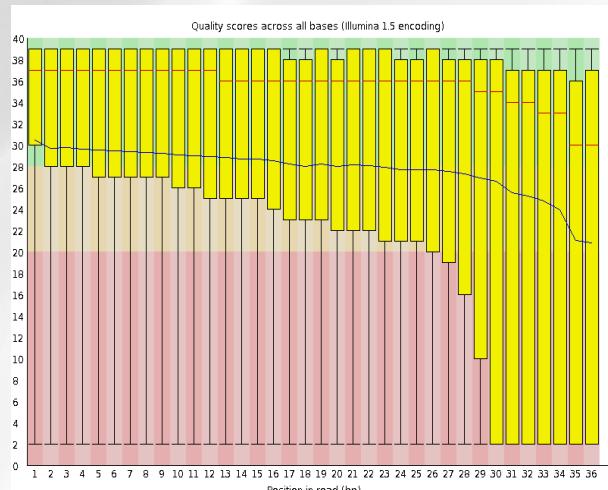
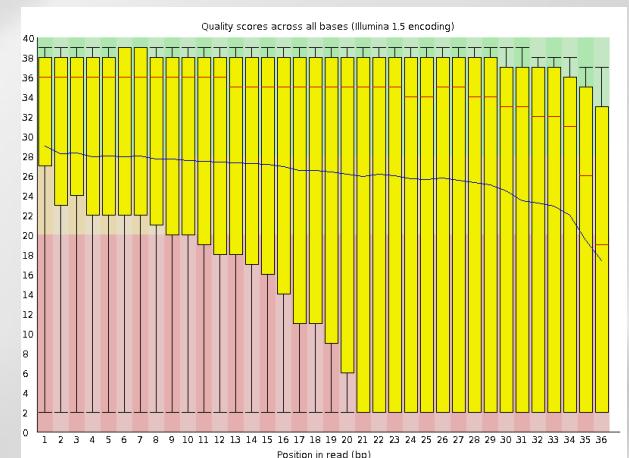
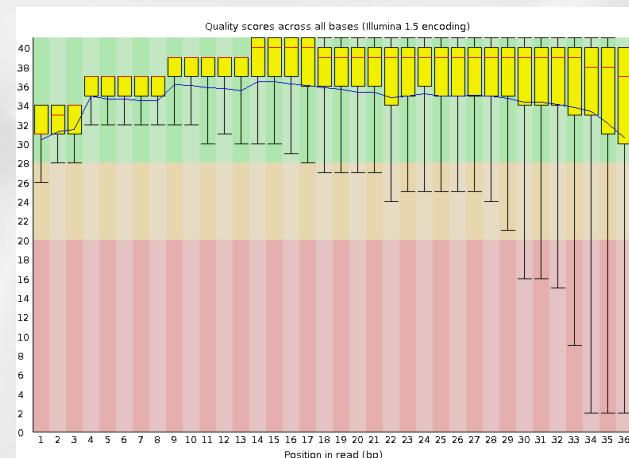
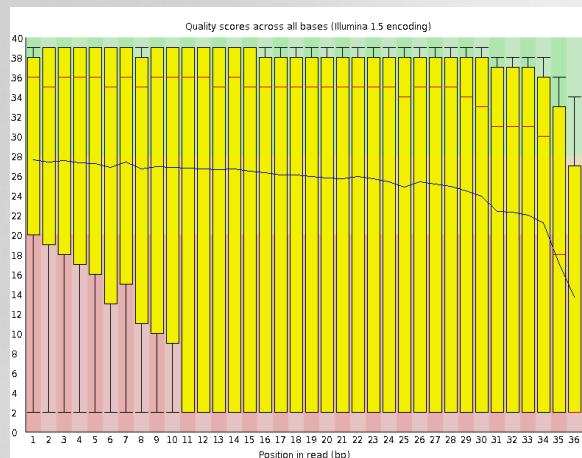
Function	A quality control tool for high throughput sequence data.
Language	Java
Requirements	A suitable Java Runtime Environment The Picard BAM/SAM Libraries (included in download)
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

A simple way to do quality control. It provides a modular set of analyses to give a quick impression of whether data has any problems of which you should be aware before doing any further analysis. The main functions of FastQC are:

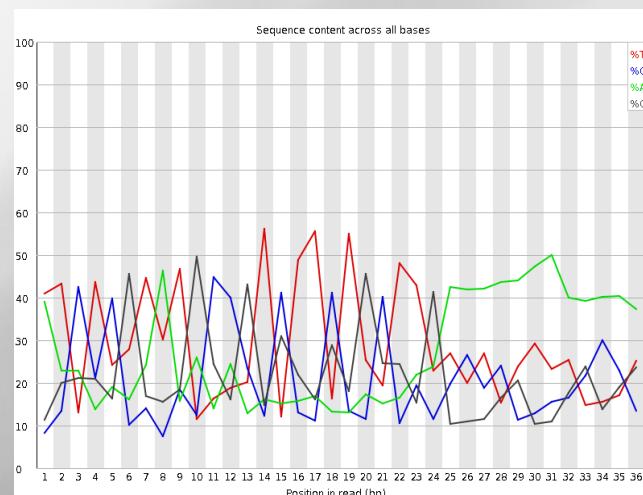
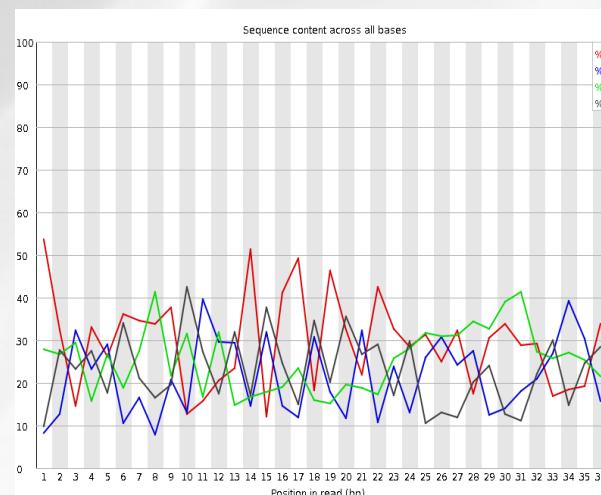
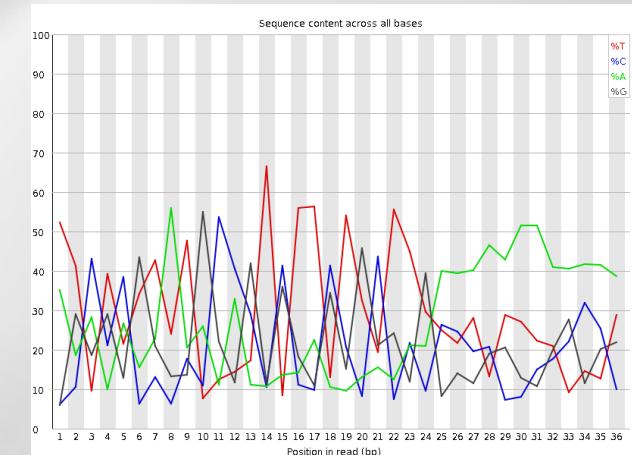
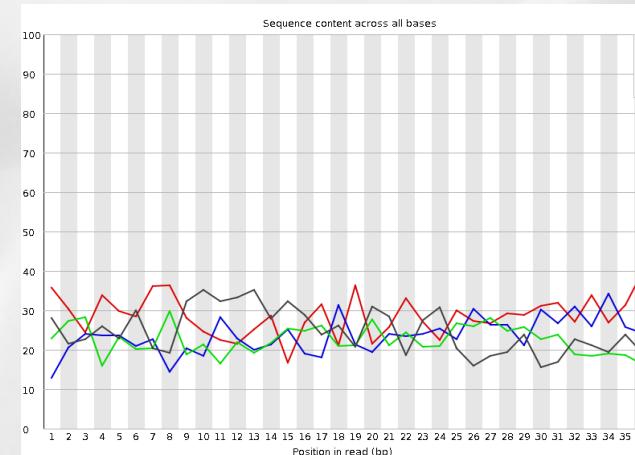
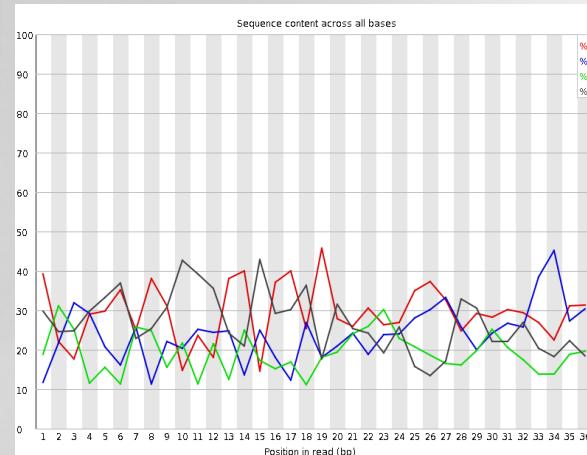
- Import of data from BAM, SAM or FastQ files (any variant)
- Provide a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive application

```
Fastqc -o nf.out nf_in.fastq
```

- Per base quality

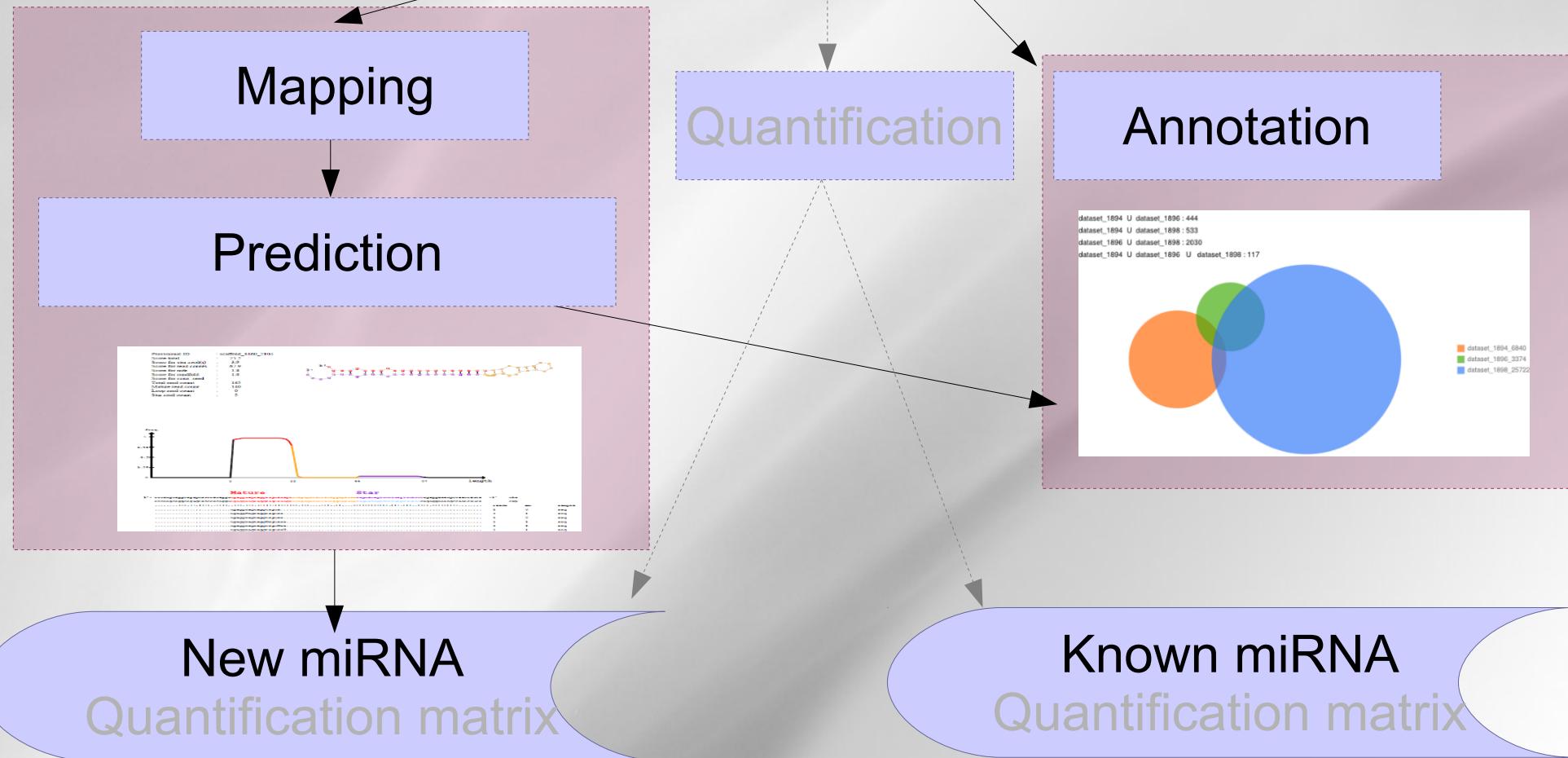


• Sequences content in nucleotides



small RNAseq pipeline

with reference



Why cleaning ?

Output reads

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 - piRNA ou tRNA
GCATTGGTGGTTCACTGGTAGAATTCTGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 - piRNA ou tRNA
GCATTGGTGGTTCACTGGTAGAATTCTGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters
- Some sequences contain sequences of interest flanked by the beginning of adapters:
 - Some of them are miRNA (yellow).

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 -piRNA ou tRNA
GCATTGGTGGTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters
- Some sequences contain sequences of interest flanked by the beginning of adapters:
 - Some of them are miRNA (yellow).
 - Some of them are other type of ncRNAs (green).

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly -N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 -piRNA ou tRNA
GCATTGGTGGTTCACTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters
- Some sequences contain sequences of interest flanked by the beginning of adapters:
 - Some of them are miRNA (yellow).
 - Some of them are other type of ncRNAs (green).
 - Some adapters contain errors (blue).

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly - N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 - piRNA ou tRNA
GCATTGGTGGTTCACTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters
- Some sequences contain sequences of interest flanked by the beginning of adapters:
 - Some of them are miRNA (yellow).
 - Some of them are other type of ncRNAs (green).
 - Some adapters contain errors (blue).
- Some sequences contain polyN (red)

```
>Adapteur
ATCTCGTATGCCGTCTCTGCTTGAAAAAAAAAAAAA
>UT1-10 - 28S rRNA
GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40 -piRNA ou tRNA
GCATTGGTGGTTCACTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT
```

Why cleaning ?

Output reads

- Some sequences contain only adapters
- Some sequences contain sequences of interest flanked by the beginning of adapters:
 - Some of them are miRNA (yellow).
 - Some of them are other type of ncRNAs (green).
 - Some adapters contain errors (blue).
- Some sequences contain polyN (red)
- Some sequences contain other type of ncRNA (pink)

>Adapteur

ATCTCGTATGCCGTCTTCTGCTTGAAAAA

>UT1-10 - 28S rRNA

GCATGTTGTGGAGAACCTGGTGCTAAATCACTCGT

>Poly-N

NN

>UT1-40 - piRNA ou tRNA

GCATTGGTGGTTCACTGGTAGAATTCTCGCCATCTC

>UT1-2-mir21

TAGCTTATCAGACTGGTGGCATCTCGTATGCCG

>UT1-3-mir143

TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT

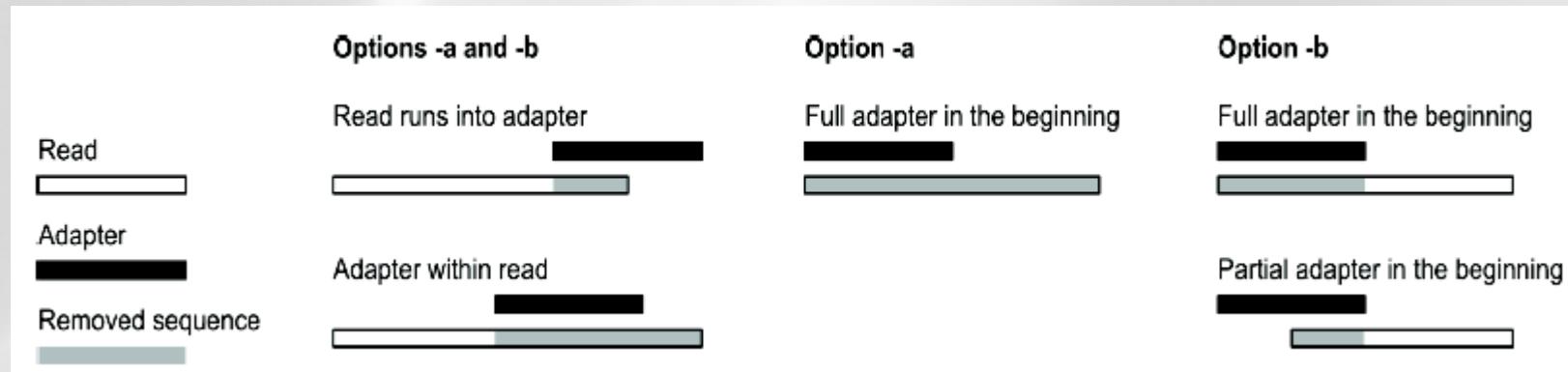
>UT1-30-mir143

TGAGATGAAGCACTGTAGCTCTCGTATGCCGTCT

• Adapters removing and length filtering

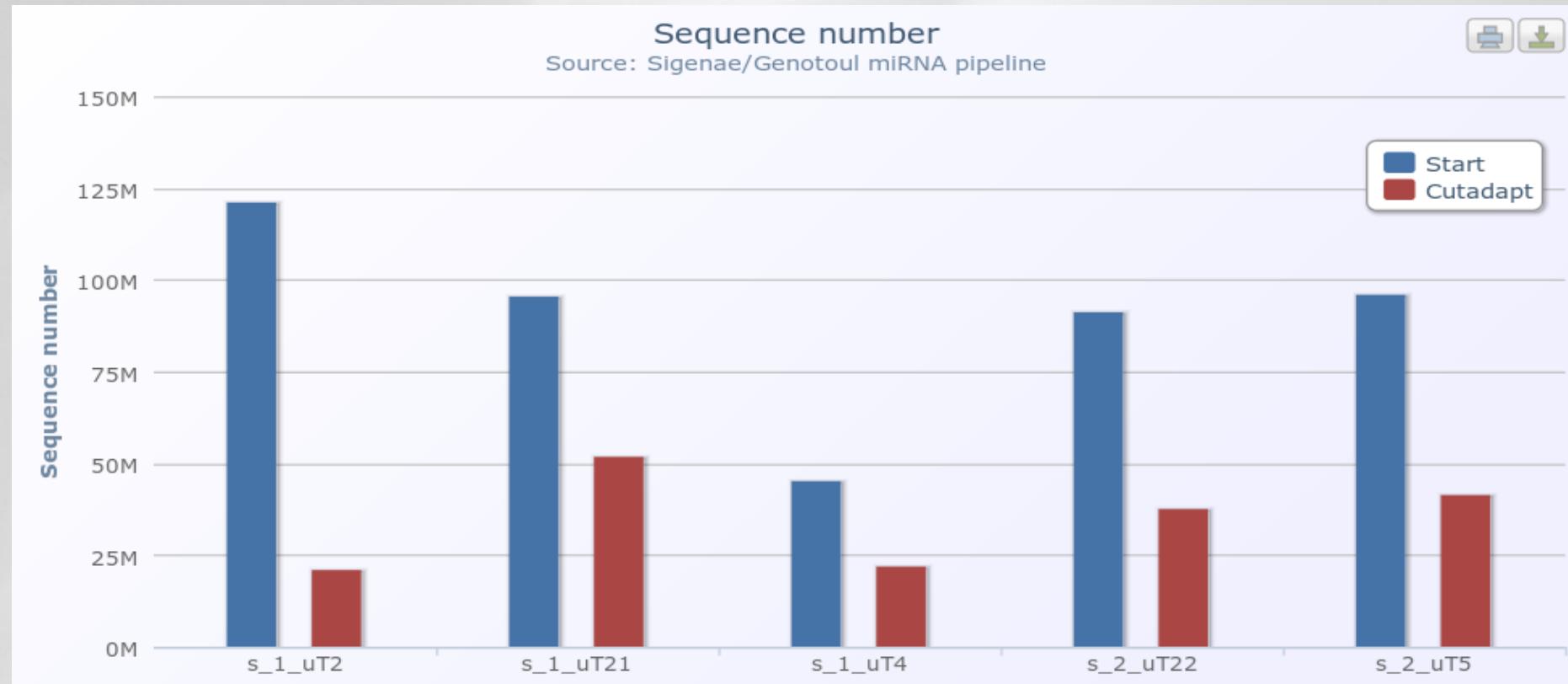
Cutadapt <http://code.google.com/p/cutadapt/>.

Cutadapt removes adapter sequences from high-throughput sequencing data. Indeed, reads are usually longer than the RNA, and therefore contain parts of the 3' adapter. It also allows to keep only sequences of desired length ($15 < \text{length} < 29$).

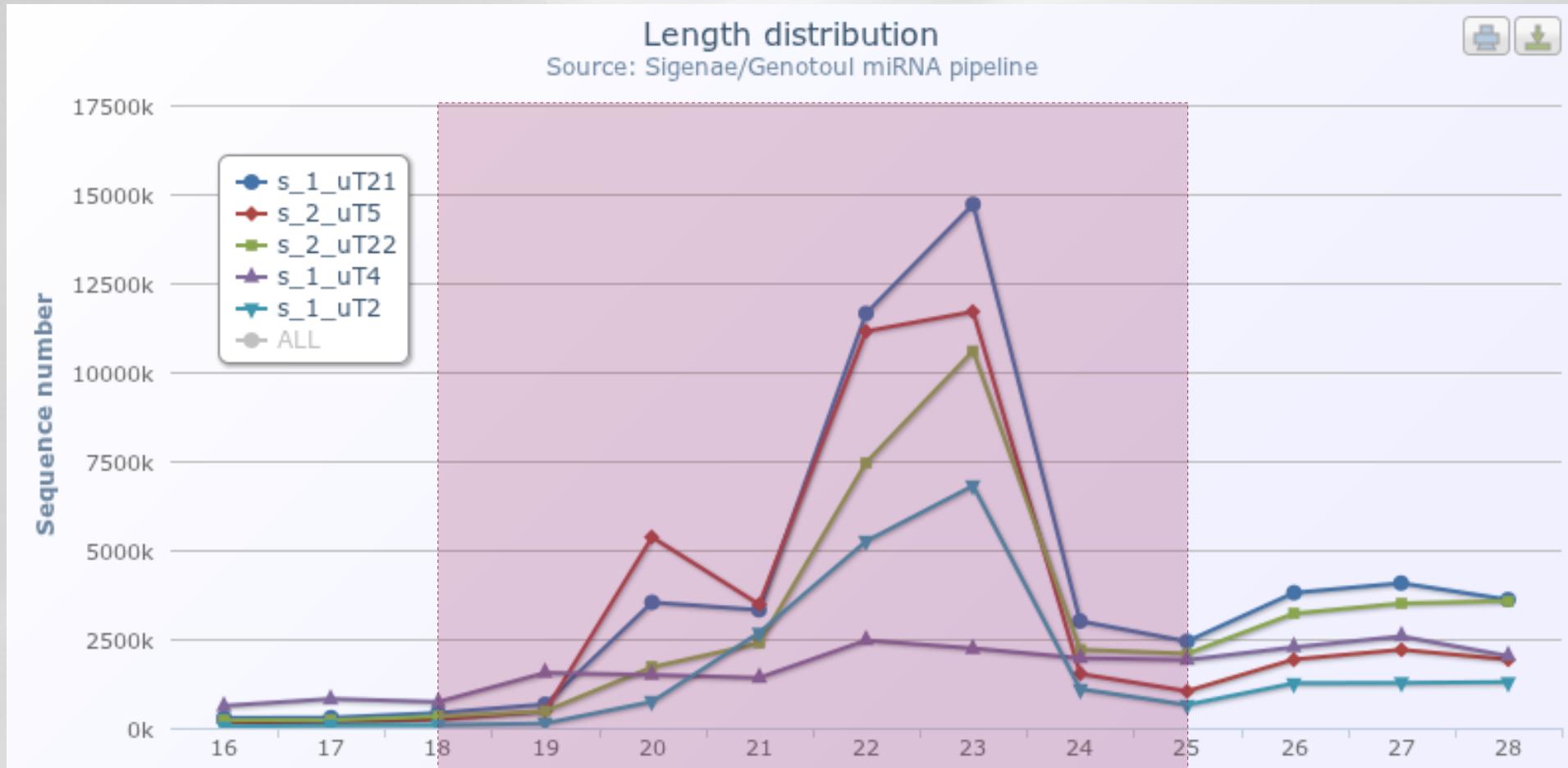


```
cutadapt -a ATCTCGTATGCCGTCTTCTGCTTG -m -M 29 -o nf_out.fq nf_in.fq
```

- 56 % of reads discarded



- Size in between 18bp:24bp
→ miRNA ?



Exercices:

- **(Re)introduction to Galaxy**
- **WF1:**
 - **Quality control**
 - **Cleaning**

Tools

Options ▾

Your user name: smaman
 Your file path : /work/smaman/

1 - UPLOAD YOUR DATA

[Get Data](#)

2 - FILES MANIPULATION

[Text Manipulation](#)[Filter and Sort](#)[Join, Subtract and Group](#)[Convert Formats](#)

3 - SEQUENCES

MANIPULATION

[FASTA manipulation](#)[FASTQ manipulation](#)[SAM/BAM manipulation : Picard \(beta\)](#)[SAM/BAM manipulation : SAM Tools](#)

4 - MAPPING

[BWA - Bowtie](#)

5 - INDEL ET SNP

[Indel Analysis](#)[RNA-Seq](#)[GATK Tools \(beta\)](#)

6 - SRNASEQ

[Analyse des miRNA](#)[Annotations](#)[Alignement sur reference](#)

WELCOME ON SIGENAE GALAXY WORKBENCH

Galaxy is a workbench available for biologists from Sigenae Platform. Galaxy objectives are:

- Make bioinfo Linux tools accessible to biologists.
 - Hide the complexity of the infrastructure.
- Allow creation, execution and sharing of workflows.

History

Options ▾

TP FastQC 54.0 Mb

[8: FastQC data 5.html](#)

[6: GM.fastqsanger](#)

[5: h1.fastqsanger](#)

[4: FastQC data 18.html](#)

[3: FASTQ Summary](#)
[Statistics on data 18](#)

[2: FASTQ Summary](#)
[Statistics on data 18](#)

76 lines, 1 comments
 format: tabular, database: ?
 Info: 99115 fastq reads were processed.

Based upon quality values and sequence characters, the input data is valid for: sanger
 Input ASCII range: '#'(35) - 'C'(67)

Input decimal range: 2 - 34

Epilog : job finished at ven mai 11 10:36:43 CEST 2012



1	2	3	4	5	6
#column	count	min	max	sum	mean
1	99115	2	33	3194703	32.2
2	99115	2	34	3156652	31.8
3	99115	2	34	3145060	31.7
4	99115	2	34	3120431	31.4
5	99115	2	34	3095075	31.3

Vos traitements bioinformatiques avec GALAXY

GALAXY pour vos traitements bioinformatiques
<http://sigenae-workbench.toulouse.inra.fr>

Une « Galaxy » parmi tant d'autres



Serveur public (<https://main.g2.bx.psu.edu/>):



- Gratuit
- Quota limité : pour se familiariser à l'outil sur des petits jeux de données.
- Données non protégées

Nombreuses autres instances :

- Curie (Nebula), URGI



Une communauté nationale et internationale très active :

- Listes de diffusion (US, FR)
- Wiki
- Twitter
- "Galaxy tour de France"

Une « Galaxy » parmi tant d'autres



L'instance locale Sigenae de Galaxy :

- Maintenue par Sigenae.
- Intégration des outils et scripts “locaux”.

→ Présentation des particularités de l'instance Sigenae.

Exemple : Analyse en 3 clics

Analyze Data Workflow Shared Data Visualization Admin Help Using 13%

Welcome smaman, you are working in /work/smaman

User Options ▾

Your user name: smaman
Your file path : /work/smaman/

1 - UPLOAD YOUR DATA

[Get Data](#)

2 - FILES MANIPULATION

[Text Manipulation](#)
[Filter and Sort](#)
[Join, Subtract and Group](#)
[Convert Formats](#)

3 - SEQUENCES
MANIPULATION

[FASTA manipulation](#)
[FASTQ manipulation](#)
[SAM/BAM manipulation : Picard
\(beta\)](#)
[SAM/BAM manipulation : SAM
Tools](#)

4 - MAPPING

[BWA - Bowtie](#)

5 - INDEL ET SNP

[Indel Analysis](#)
[RNA-Seq](#)
[GATK Tools \(beta\)](#)

History Options ▾

Unnamed history 0 bytes

Your history is empty. Click 'Get Data' on the left pane to start

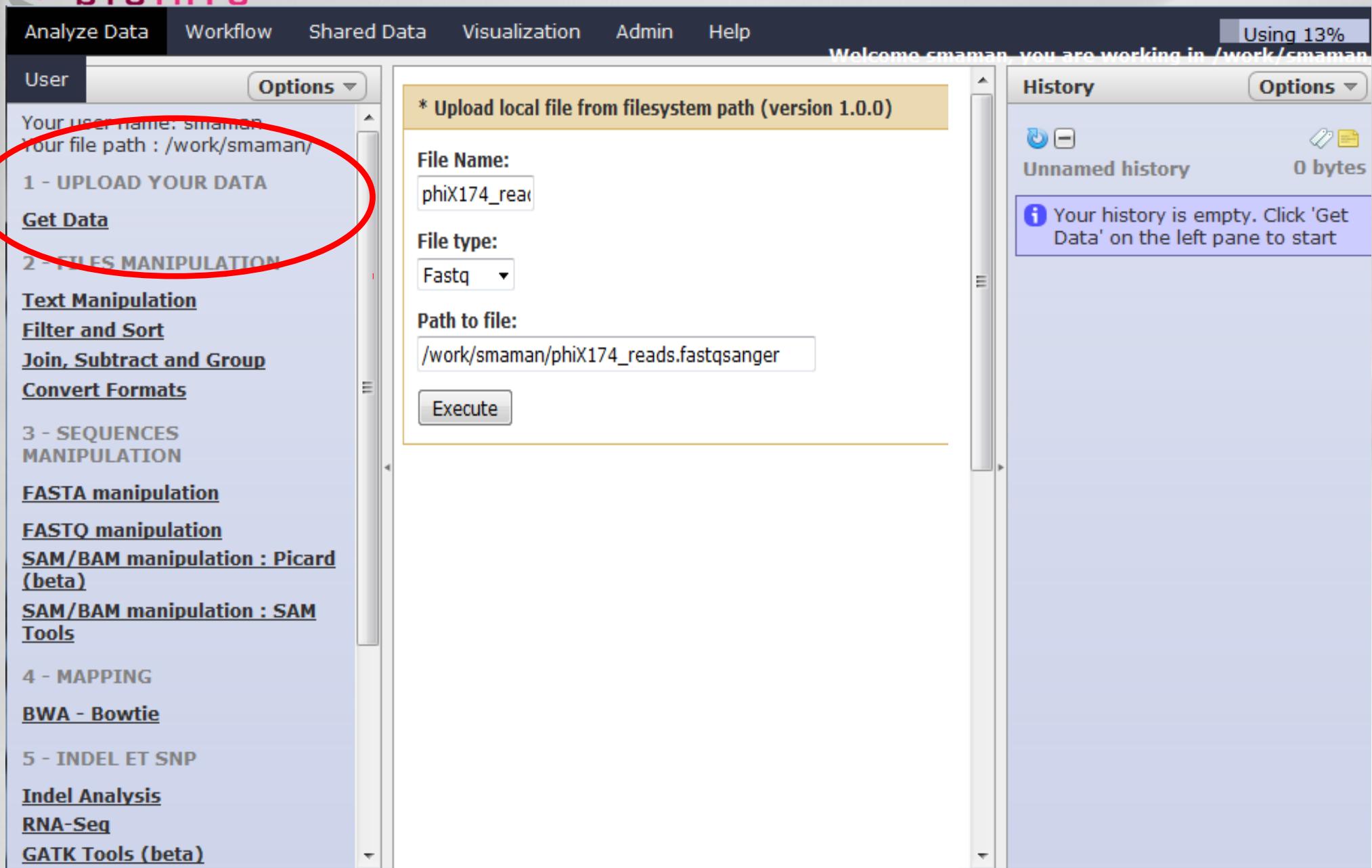
WELCOME ON SIGENAE GALAXY WORKBENCH

Galaxy is a workbench available for biologists from Sigenae Platform. Galaxy objectives are:

- Make bioinfo Linux tools accessible to biologists.
 - Hide the complexity of the infrastructure.
 - Allow creation, execution and sharing of workflows.



Exemple : Analyse en 3 clics



The screenshot shows the genotoul bioinfo software interface. The top navigation bar includes 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Admin', and 'Help'. A status bar at the top right indicates 'Using 13%' and 'Welcome smaman, you are working in /work/smaman'. The main area has three panes:

- Left Sidebar (User):** Shows the user name 'smaman' and file path '/work/smaman/'. It lists categories: 1 - UPLOAD YOUR DATA (with a red circle around it), 2 - FILES MANIPULATION, 3 - SEQUENCES MANIPULATION, 4 - MAPPING, 5 - INDEL ET SNP, and RNA-Seq/GATK Tools (beta).
- Middle Panel:** Titled '* Upload local file from filesystem path (version 1.0.0)'. It contains fields for 'File Name' (phiX174_reads), 'File type' (Fastq), and 'Path to file' (/work/smaman/phiX174_reads.fastqsanger). An 'Execute' button is at the bottom.
- Right Panel:** Titled 'History'. It shows an 'Unnamed history' entry with 0 bytes. A message says 'Your history is empty. Click 'Get Data' on the left pane to start'.

Exemple : Analyse en 3 clics

Analyze Data Workflow Shared Data Visualization Admin Help

Welcome smaman, you are working in /work/smaman
Using 13%

User Options

Your user name: smaman
Your file path : /work/smaman/

1 - UPLOAD YOUR DATA

[Get Data](#)

2 - FILES MANIPULATION

[Text Manipulation](#)
[Filter and Sort](#)
[Join, Subtract and Group](#)
[Convert Formats](#)

3 - SEQUENCES MANIPULATION

[FASTA manipulation](#)
[FASTQ manipulation](#)
[SAM/BAM manipulation : Picard \(beta\)](#)
[SAM/BAM manipulation : SAM Tools](#)

4 - MAPPING

[BWA - Bowtie](#)

5 - INDEL ET SNP

[Indel Analysis](#)
[RNA-Seq](#)
[GATK Tools \(beta\)](#)

* Upload local file from filesystem path (version 1.0.0)

File Name: phiX174_reads

File type: Fastq

Path to file: /work/smaman/phiX174_reads.fastqsanger

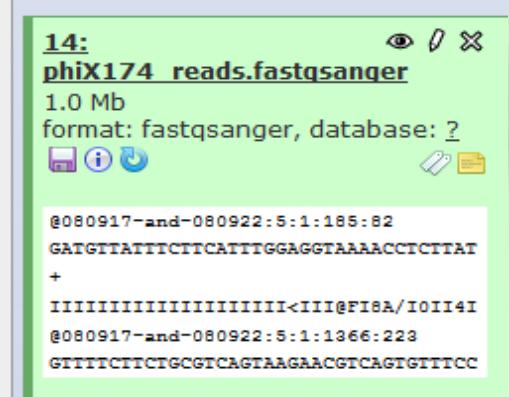
Execute

History Options

Unnamed history 0 bytes

Your history is empty. Click 'Get Data' on the left pane to start

Galaxy sensibilisation - TP 12.1 Mb
2 - BWA and FastQC

14: phiX174_reads.fastqsanger 1.0 Mb
format: fastqsanger, database: ?


```

@080917-and-080922:5:1:185:82
GATGTTATTCTTCATTGGAGGTAAACCTCTTAT
+
IIIIIIIIIIIIIIII<III@F16A/I0II4I
@080917-and-080922:5:1:1366:223
GTTTCTCTGCGTCAGTAAGAACGTCAGTGTTC

```

Exemple : Analyse en 3 clics

Analyze Data Workflow Shared Data Visualization Admin Help Using 13%

Welcome smaman, you are working in /work/smaman

User Options

Your user name: smaman
Your file path : /work/smaman/

1 - UPLOAD YOUR DATA
[Get Data](#)

2 - FILES MANIPULATION

[Text Manipulation](#)
[Filter and Sort](#)
[Join, Subtract and Group](#)

[NGS: Mapping](#)

- [Lastz map short reads against reference sequence](#)
- [Lastz paired reads map short paired reads against reference sequence](#)
- [Map with Bowtie for Illumina](#)
- [Map with Bowtie for SOLiD](#)
- [Map with BWA for Illumina](#)

4 - MAPPING

[BWA - Bowtie](#)

5 - INDEL ET SNP

[Indel Analysis](#)
[RNA-Seq](#)
[GATK Tools \(beta\)](#)

History Options

Unnamed history 0 bytes

Your history is empty. Click 'Get Data' on the left pane to start

WELCOME ON SIGENAE GALAXY WORKBENCH

Galaxy is a workbench available for biologists from Sigenae Platform. Galaxy objectives are:

- Make bioinfo Linux tools accessible to biologists.
 - Hide the complexity of the infrastructure.
 - Allow creation, execution and sharing of workflows.

Exemple : Analyse en 3 clics

Analyze Data Workflow Shared Data Visualization Admin Help Using 13%

Welcome smaman, you are working in /work/smaman

User Options

Your user name: smaman
Your file path : /work/smaman/

1 - UPLOAD YOUR DATA
[Get Data](#)

2 - FILES MANIPULATION
[Text Manipulation](#)
[Filter and Sort](#)
[Join, Subtract and Group](#)

NGS: Mapping

- [Lastz map short reads against reference sequence](#)
- [Lastz paired reads map short paired reads against reference sequence](#)
- [Map with Bowtie for Illumina](#)
- [Map with Bowtie for SOLiD](#)
- [Map with BWA for Illumina](#)

4 - MAPPING
[BWA - Bowtie](#)

5 - INDEL ET SNP
[Indel Analysis](#)
[RNA-Seq](#)
[GATK Tools \(beta\)](#)

Map with BWA for Illumina (version 1.2.2)

Will you select a reference genome from your history?
Use one from the history ▾

Select a reference from history:
11: phiX174_genome.fa ▾

Is this library mate-paired?:
Single-end ▾

FASTQ file:
14: phiX174_reads.fastqsanger ▾

FASTQ with either Sanger-scaled quality values (f)

History Options

Unnamed history 0 bytes

i Your history is empty. Click 'Get Data' on the left pane to start

Exemple : Analyse en 3 clics

Analyze Data Workflow Shared Data Visualization Admin Help Using 13%

Welcome smaman, you are working in /work/smaman

User Options

Your user name: smaman
Your file path : /work/smaman/

1 - UPLOAD YOUR DATA
[Get Data](#)

2 - FILES MANIPULATION
[Text Manipulation](#)
[Filter and Sort](#)
[Join, Subtract and Group](#)
[Convert Formats](#)

3 - SEQUENCES
MANIPULATION
[FASTA manipulation](#)
[FASTQ manipulation](#)
[SAM/BAM manipulation : Picard
\(beta\)](#)
[SAM/BAM manipulation : SAM
Tools](#)

4 - MAPPING
[BWA - Bowtie](#)

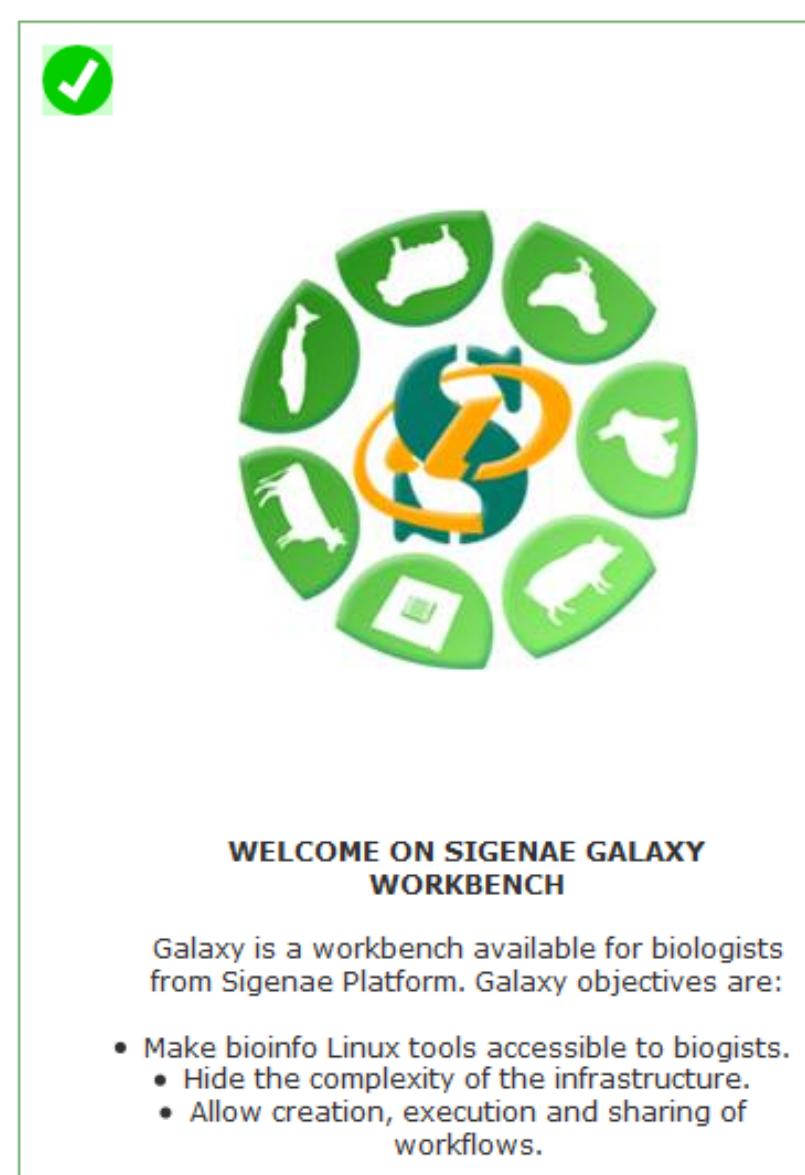
5 - INDEL ET SNP
[Indel Analysis](#)
[RNA-Seq](#)
[GATK Tools \(beta\)](#)

History Options

Unnamed history 0 bytes

Your history is empty. Click 'Get Data' on the left pane to start

15: Map with BWA for [Illumina on data 14](#) and [data 11: mapped reads](#)
Job is waiting to run



WELCOME ON SIGENAE GALAXY WORKBENCH

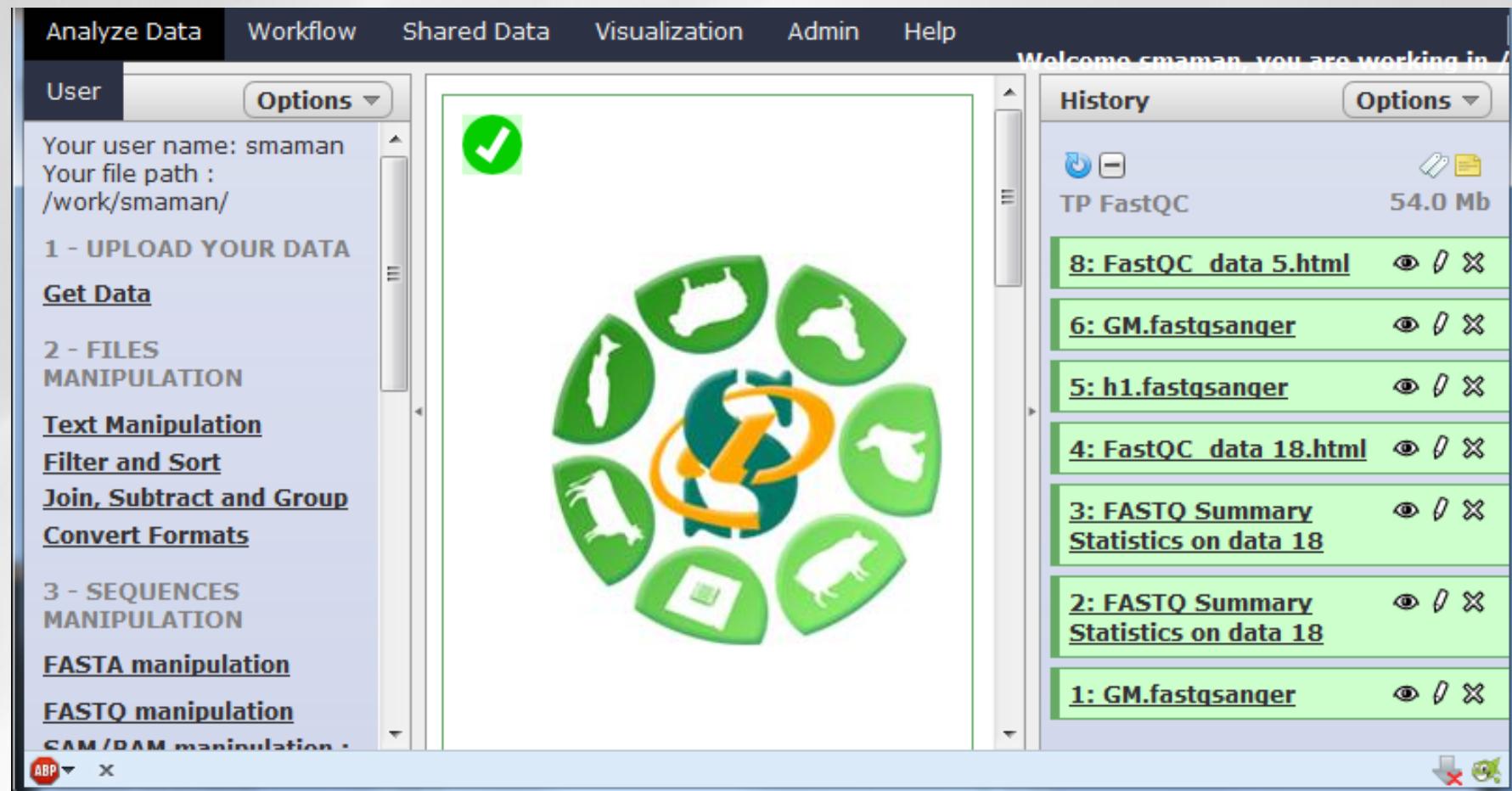
Galaxy is a workbench available for biologists from Sigenae Platform. Galaxy objectives are:

- Make bioinfo Linux tools accessible to biologists.
 - Hide the complexity of the infrastructure.
 - Allow creation, execution and sharing of workflows.

Interface intuitive

Interface divisée en 4 parties :

- 1 - Liste des outils disponibles.
- 2 - Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 - Historique ou workflow détaillé.
- 4 - Menu .



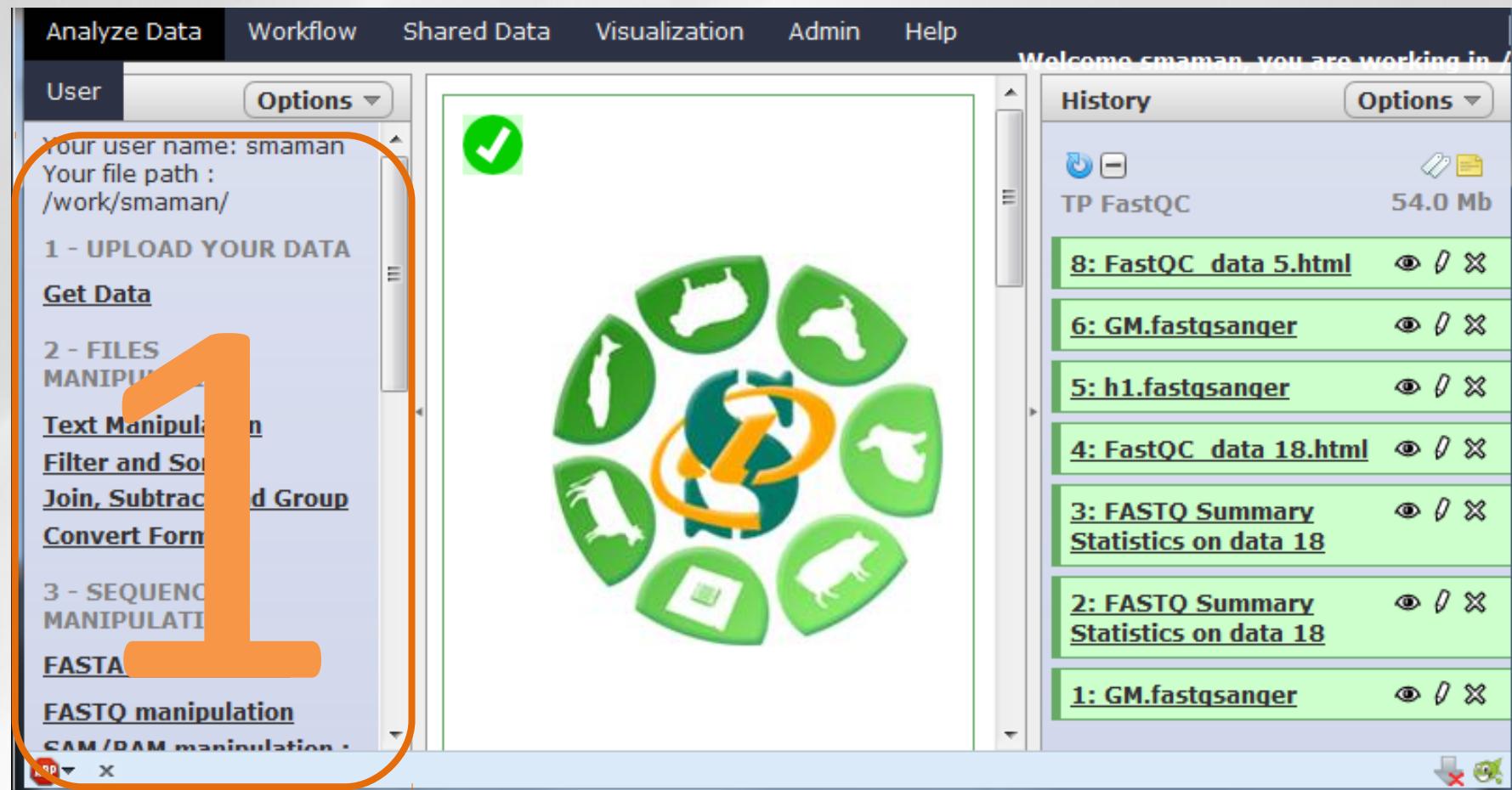
The screenshot shows the software interface divided into four main sections:

- Analyze Data Panel (Left):** Contains a user profile section with the name "smaman" and file path "/work/smaman/". It lists several categories of tools:
 - 1 - UPLOAD YOUR DATA**: Includes links for "Get Data", "Text Manipulation", "Filter and Sort", "Join, Subtract and Group", and "Convert Formats".
 - 2 - FILES MANIPULATION**: Includes links for "Text Manipulation", "Filter and Sort", "Join, Subtract and Group", and "Convert Formats".
 - 3 - SEQUENCES MANIPULATION**: Includes links for "FASTA manipulation", "FASTQ manipulation", and "SAM/BAM manipulation".
- Workflow Panel (Center):** Displays a green checkmark icon above a circular diagram consisting of several green leaf-like shapes containing white silhouettes of various animals (cow, sheep, pig, etc.) and a central orange recycling symbol.
- Shared Data Panel (Top Right):** Shows a "History" list with the following items:
 - TP FastQC (54.0 Mb)
 - 8: FastQC data 5.html
 - 6: GM.fastqsanger
 - 5: h1.fastqsanger
 - 4: FastQC data 18.html
 - 3: FASTQ Summary Statistics on data 18
 - 2: FASTQ Summary Statistics on data 18
 - 1: GM.fastqsanger
- Help Panel (Bottom Right):** Contains standard help icons: ABP, X, and a magnifying glass.

Interface divisée en 4 parties :

- 1 - Liste des outils disponibles.
- 2 - Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 - Historique ou workflow détaillé.
- 4 - Menu .

1

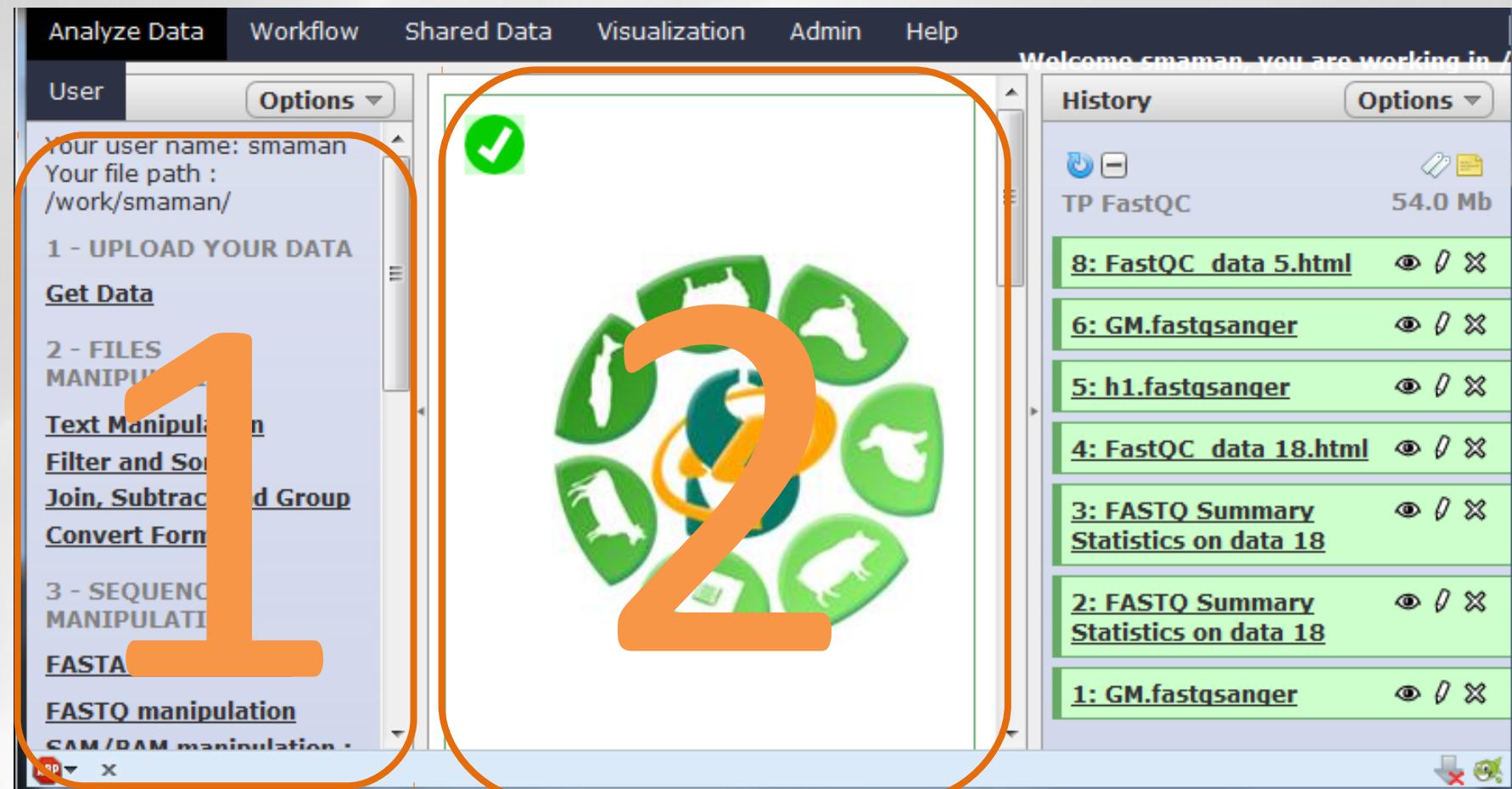


Step	Description
TP FastQC	54.0 Mb
8: FastQC data 5.html	View Edit Delete
6: GM.fastqsanger	View Edit Delete
5: h1.fastqsanger	View Edit Delete
4: FastQC data 18.html	View Edit Delete
3: FASTQ Summary Statistics on data 18	View Edit Delete
2: FASTQ Summary Statistics on data 18	View Edit Delete
1: GM.fastqsanger	View Edit Delete

Interface intuitive

Interface divisée en 4 parties :

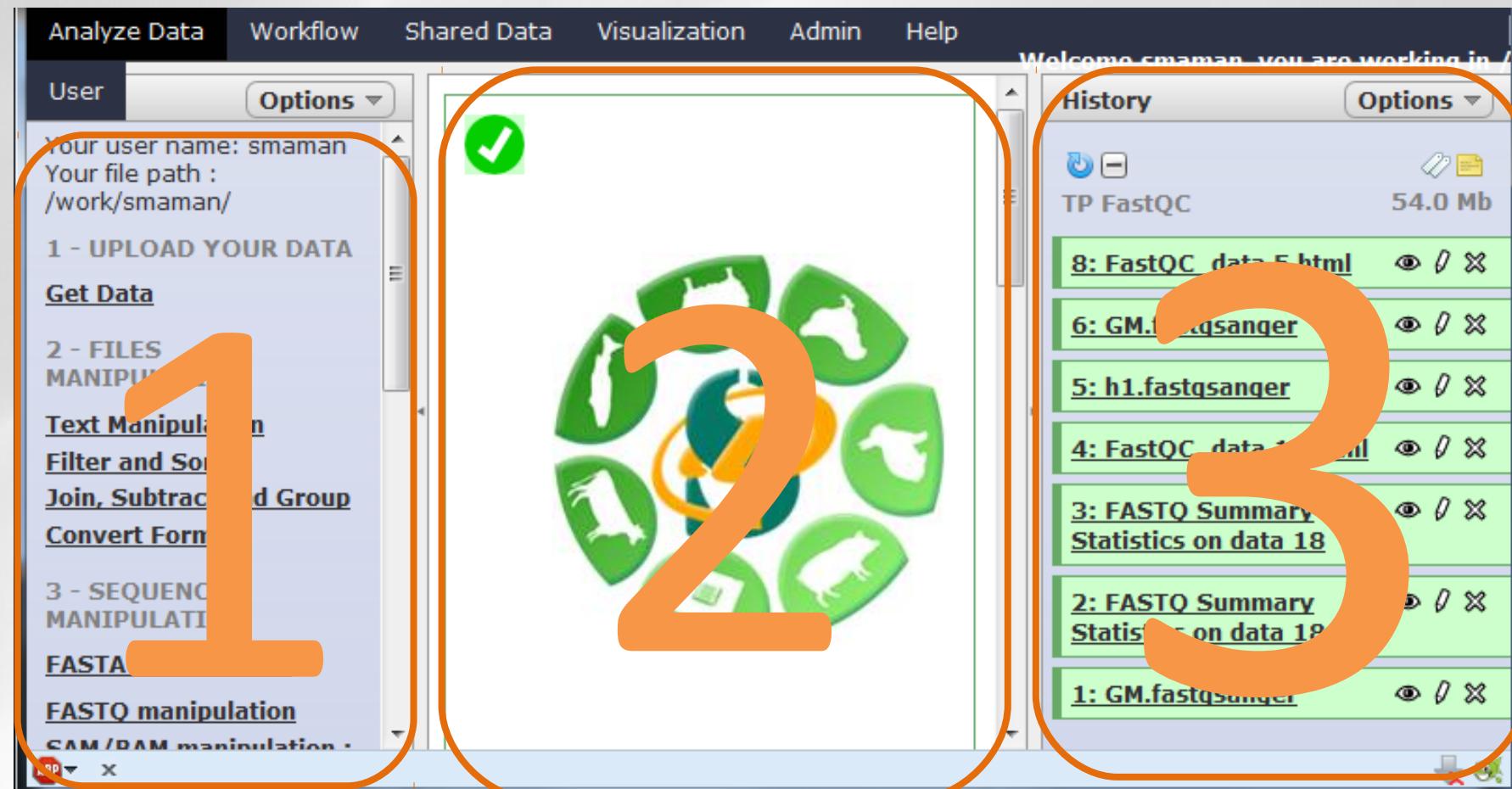
- 1 - Liste des outils disponibles.
- 2 - Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 - Historique ou workflow détaillé.
- 4 - Menu .



Interface intuitive

Interface divisée en 4 parties :

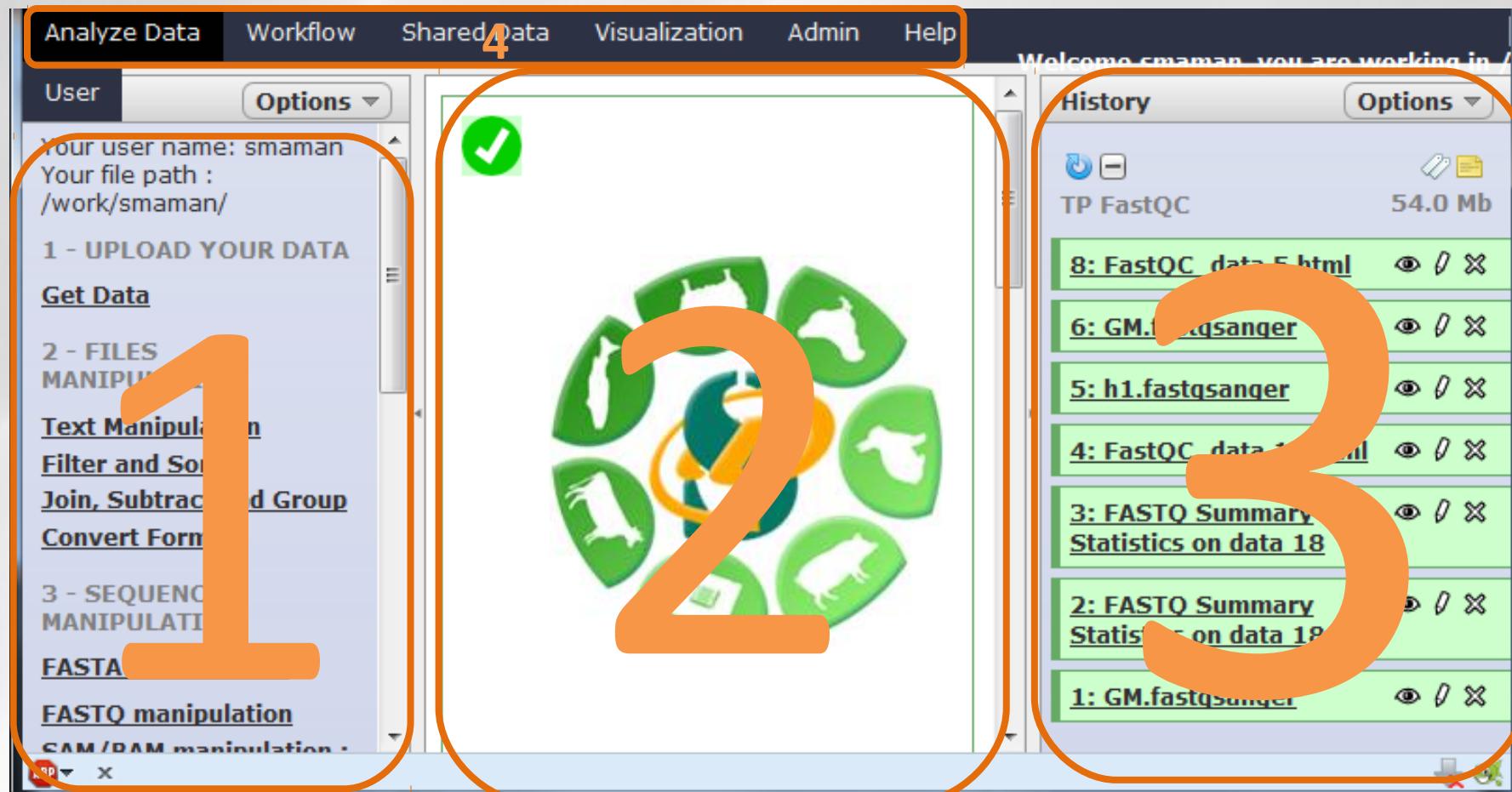
- 1 - Liste des outils disponibles.
- 2 - Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 - Historique ou workflow détaillé.
- 4 - Menu .



Interface intuitive

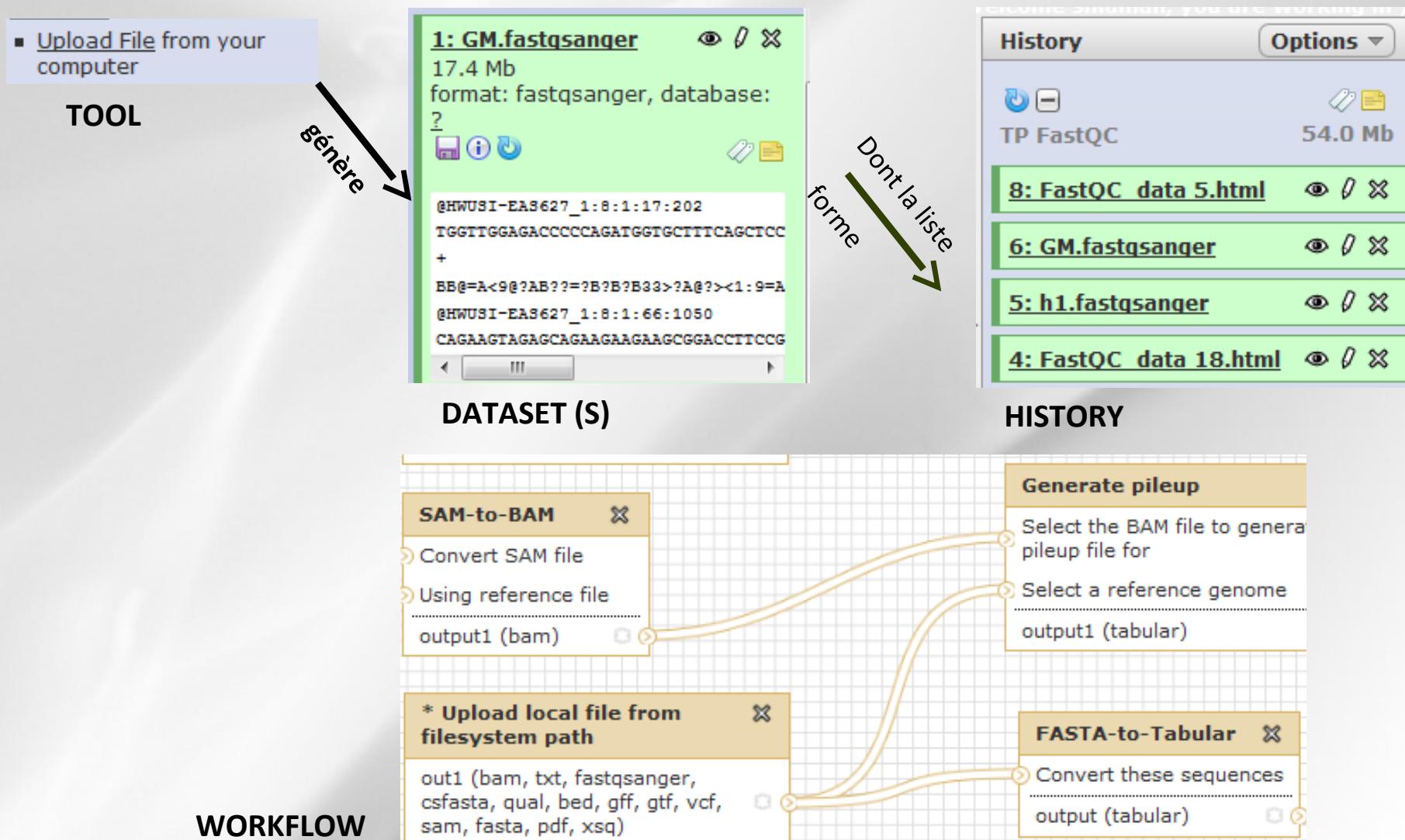
Interface divisée en 4 parties :

- 1 - Liste des outils disponibles.
- 2 - Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 - Historique ou workflow détaillé.
- 4 - Menu .



Vocabulaire spécifique à Galaxy

- TOOL** : Outil bioinformatique ou de traitement de fichiers.
- DATASET** : Fichier téléchargé dans Galaxy (entrant) ou fichier généré (résultat).
- HISTORY** : Liste des datasets (entrants et résultants) générés par les tools.
- WORKFLOW** : Chaîne de traitements (visualisation, édition, partage, etc.)



Présentation générale des workflows

Gestion du workflow

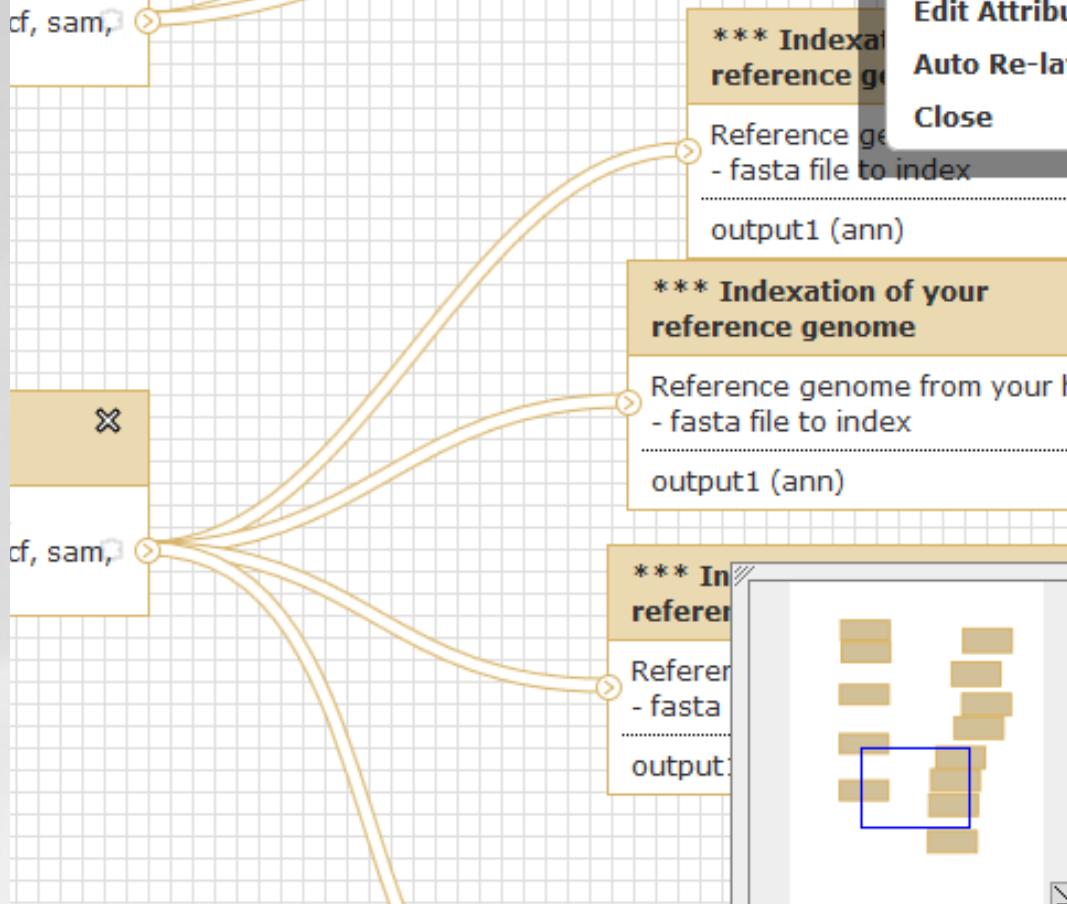
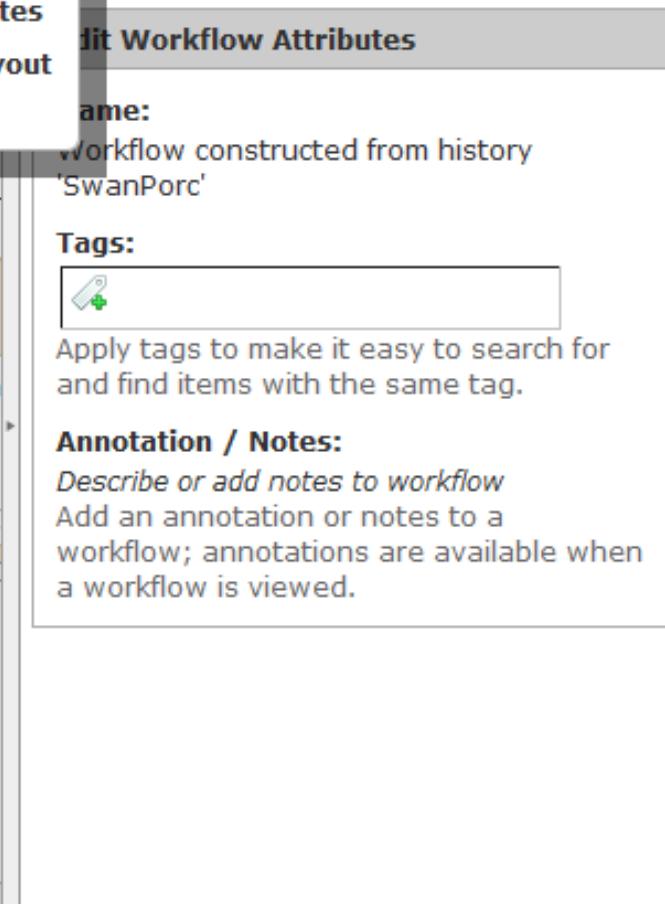
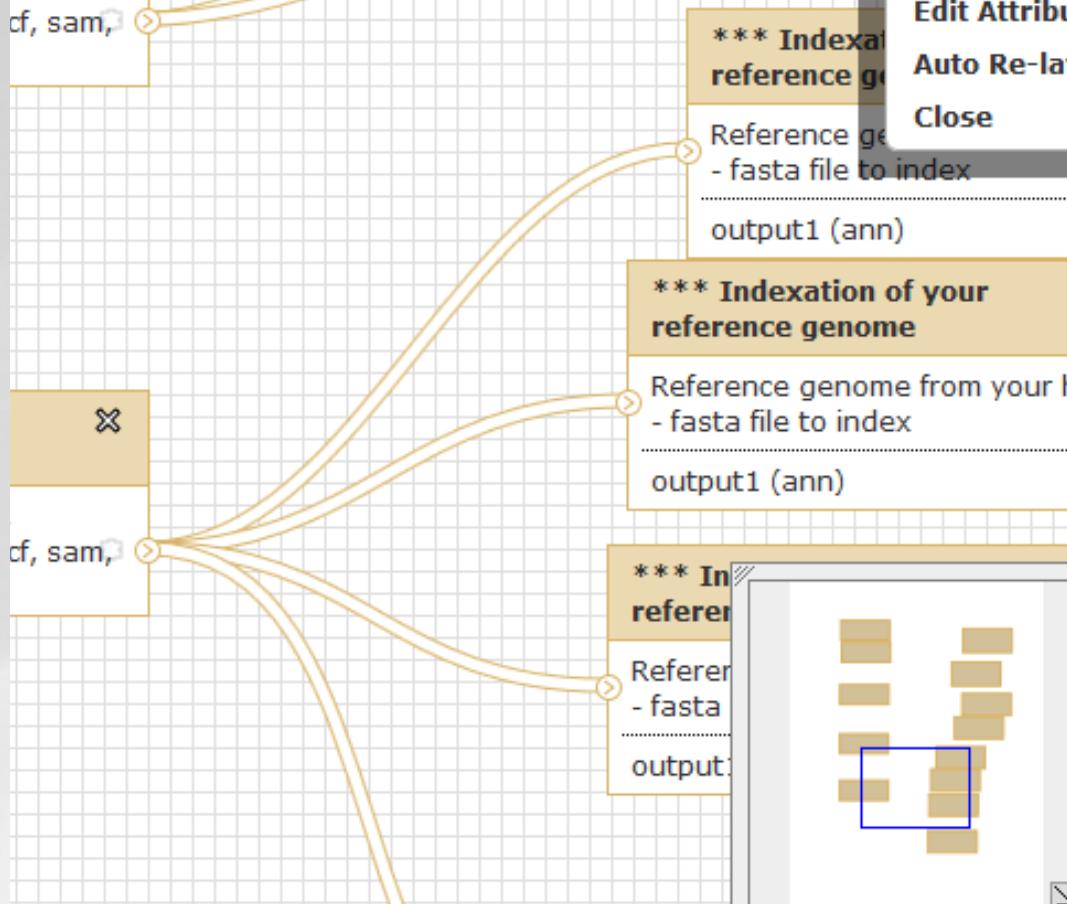


Diagramme de flux



Zoom



Ajout d'attributs

Canvas | Workflow constructed from history 'SwanPorc'

Save **Details**

Edit Workflow Attributes

Name: Workflow constructed from history 'SwanPorc'

Tags:

Apply tags to make it easy to search for and find items with the same tag.

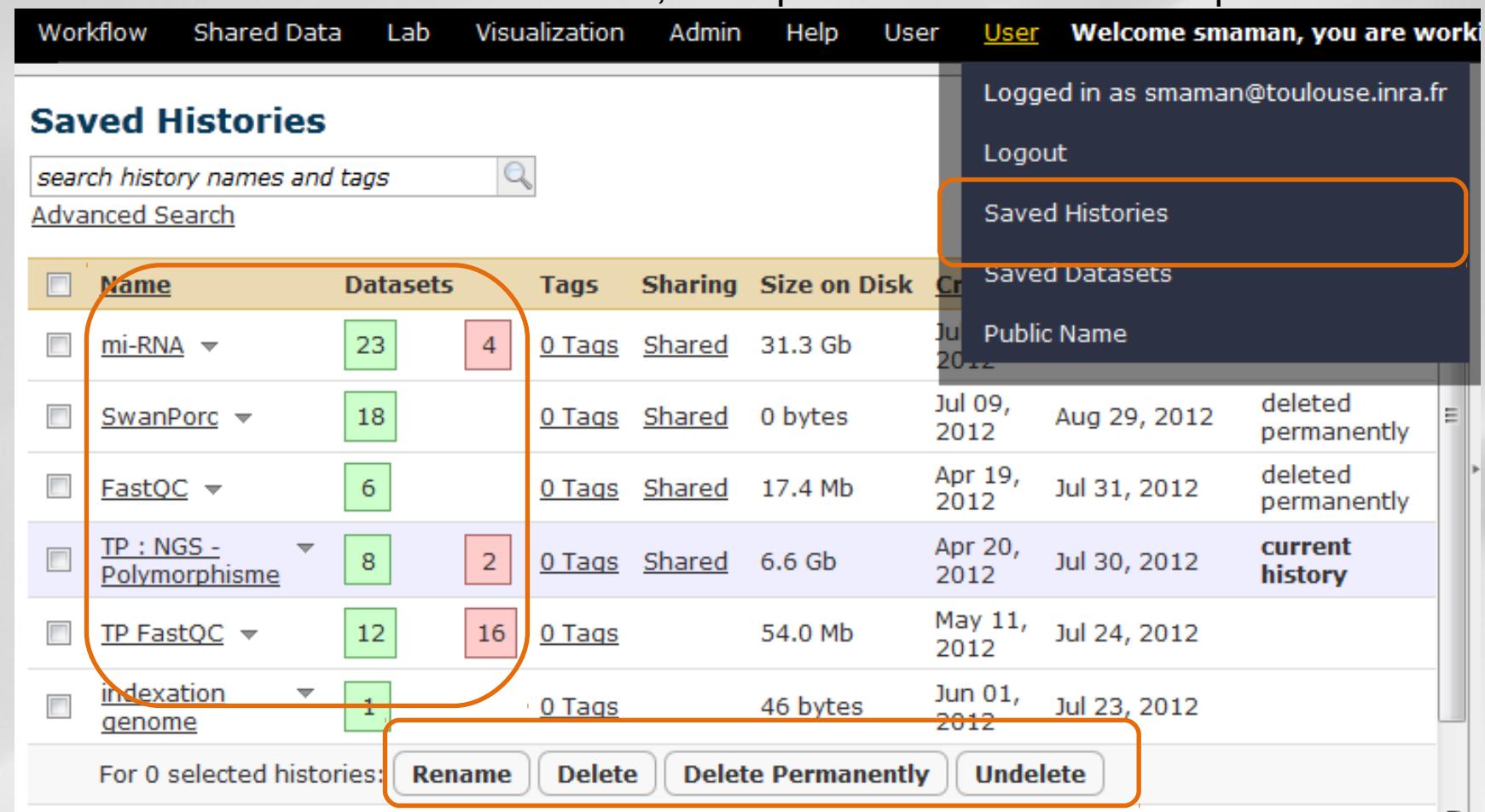
Annotation / Notes:

Describe or add notes to workflow

Add an annotation or notes to a workflow; annotations are available when a workflow is viewed.

Gestion de vos historiques

Depuis le menu « User » / « Saved Histories », vous pouvez lister vos historiques :



Name	Datasets	Tags	Sharing	Size on Disk	Creation Date	Last Update	
mi-RNA ▾	23	4	0 Tags Shared	31.3 Gb	Jul 01, 2012	Aug 29, 2012	Public Name
SwanPorc ▾	18		0 Tags Shared	0 bytes	Jul 09, 2012	Aug 29, 2012	deleted permanently
FastQC ▾	6		0 Tags Shared	17.4 Mb	Apr 19, 2012	Jul 31, 2012	deleted permanently
TP : NGS - Polymorphisme ▾	8	2	0 Tags Shared	6.6 Gb	Apr 20, 2012	Jul 30, 2012	current history
TP FastQC ▾	12	16	0 Tags	54.0 Mb	May 11, 2012	Jul 24, 2012	
indexation genome ▾	1		0 Tags	46 bytes	Jun 01, 2012	Jul 23, 2012	

For 0 selected histories:

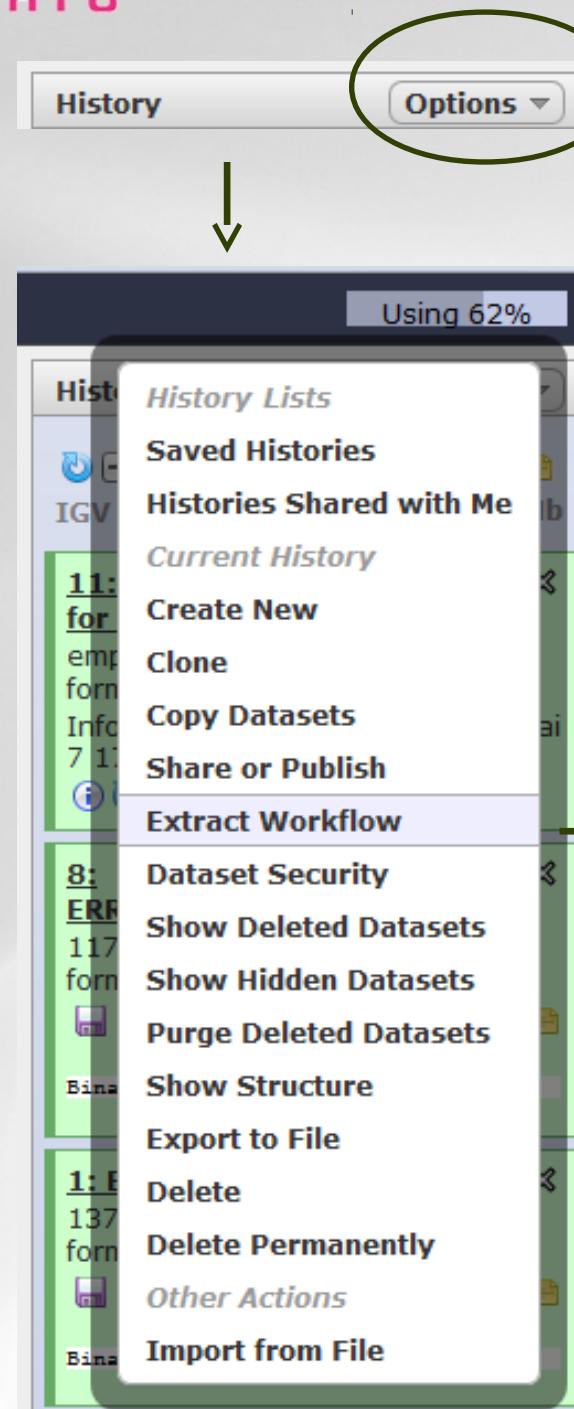
Rename **Delete** **Delete Permanently** **Undelete**



Un simple clic droit sur le titre de votre historique vous permet de le gérer : renommer, supprimer (définitivement ou pas) et rétablir.

Pour travailler sur un historique, sélectionner l'option « View ».

Exporter vos historiques en workflows



Depuis votre fenêtre « History » :

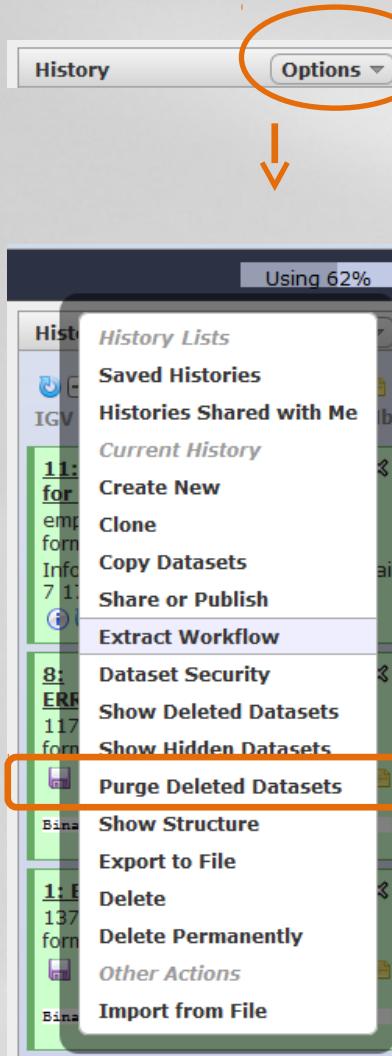
- 1 – Cliquer sur « Options » .
- 2 – « Extract workflow » .
- 3 – Nommer votre workflow.
- 4 – Sélectionner les tools utiles.
- 5 – « Create workflow » .

Workflow name
Workflow constructed from history 'IGV bai'

Tool

	History items created
* Upload local file from filesystem path <input checked="" type="checkbox"/> Include "* Upload local file from filesystem path" in workflow	▶ 1: ERR000017.bam
* Upload local file from filesystem path <input checked="" type="checkbox"/> Include "* Upload local file from filesystem path" in workflow	▶ 8: ERR000017.sorted
* BAM sorted to BAI for IGV <input checked="" type="checkbox"/> Include "* BAM sorted to BAI for IGV" in workflow	▶ 11: * BAM sorted to

Comment être un bon Galaxy user ?



ORGANISER SON ESPACE DE TRAVAIL :

Créer un nouvel historique par analyse.

Renommer et taguer vos fichiers résultats, vos historiques et vos workflows.

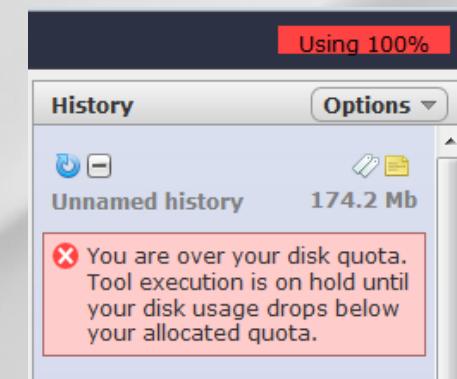
MAITRISER SON QUOTA :

Télécharger vos données privées sans copie sur le serveur.

Ne pas confondre datasets supprimés et datasets supprimés définitivement.

Supprimer définitivement et régulièrement les datasets inutiles.

Supprimer vos historiques et vos workflows inutiles.



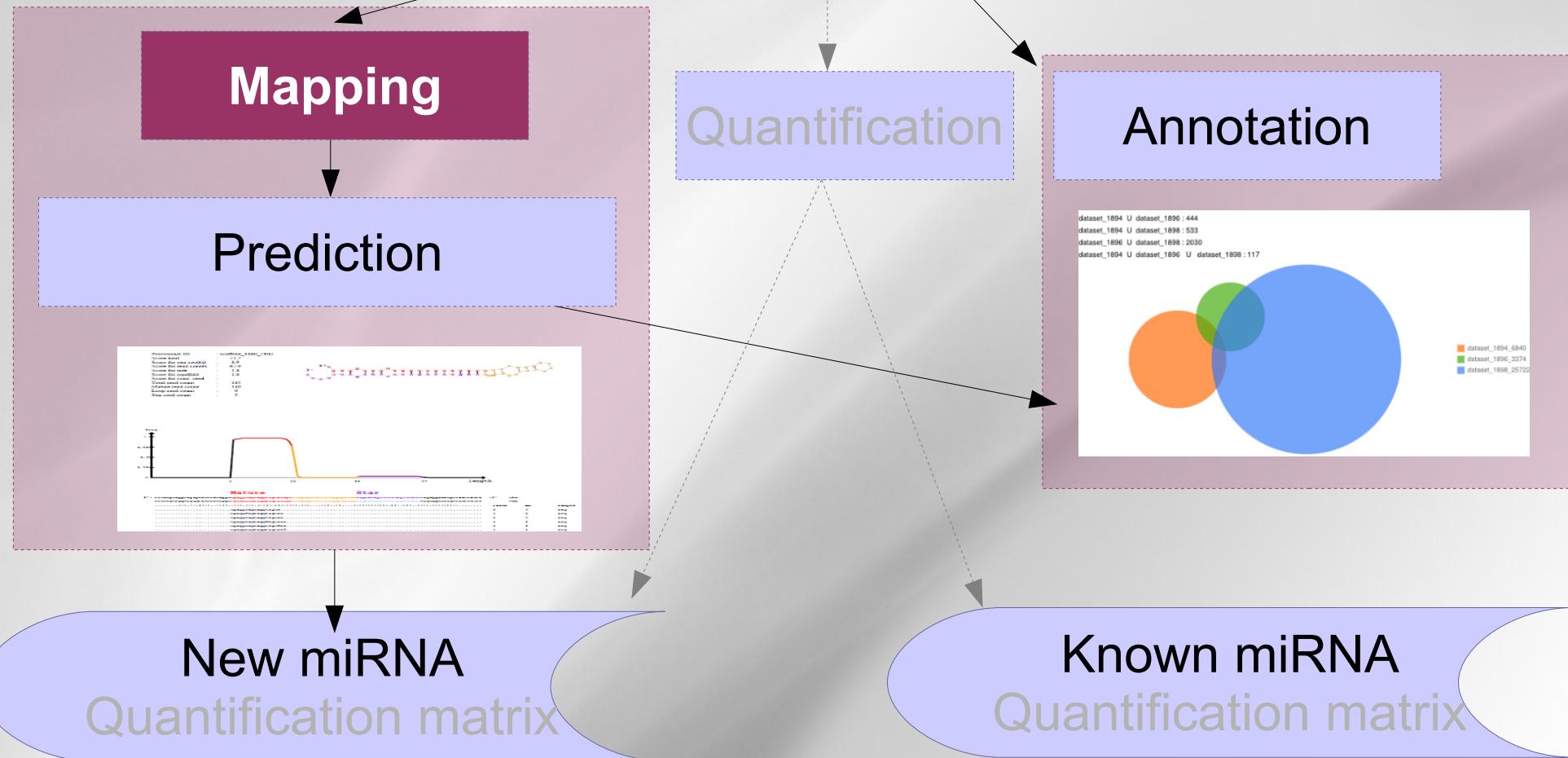
N'HESITEZ PAS A DEMANDER DE L'AIDE :

Via un ticket sur la forge DGA ou un mail (sigenae-support@listes.inra.fr).

Consulter la FAQ, le manuel utilisateur et les supports de formation sur « sig-learning ».

small RNAseq pipeline

with reference



New miRNA
Quantification matrix

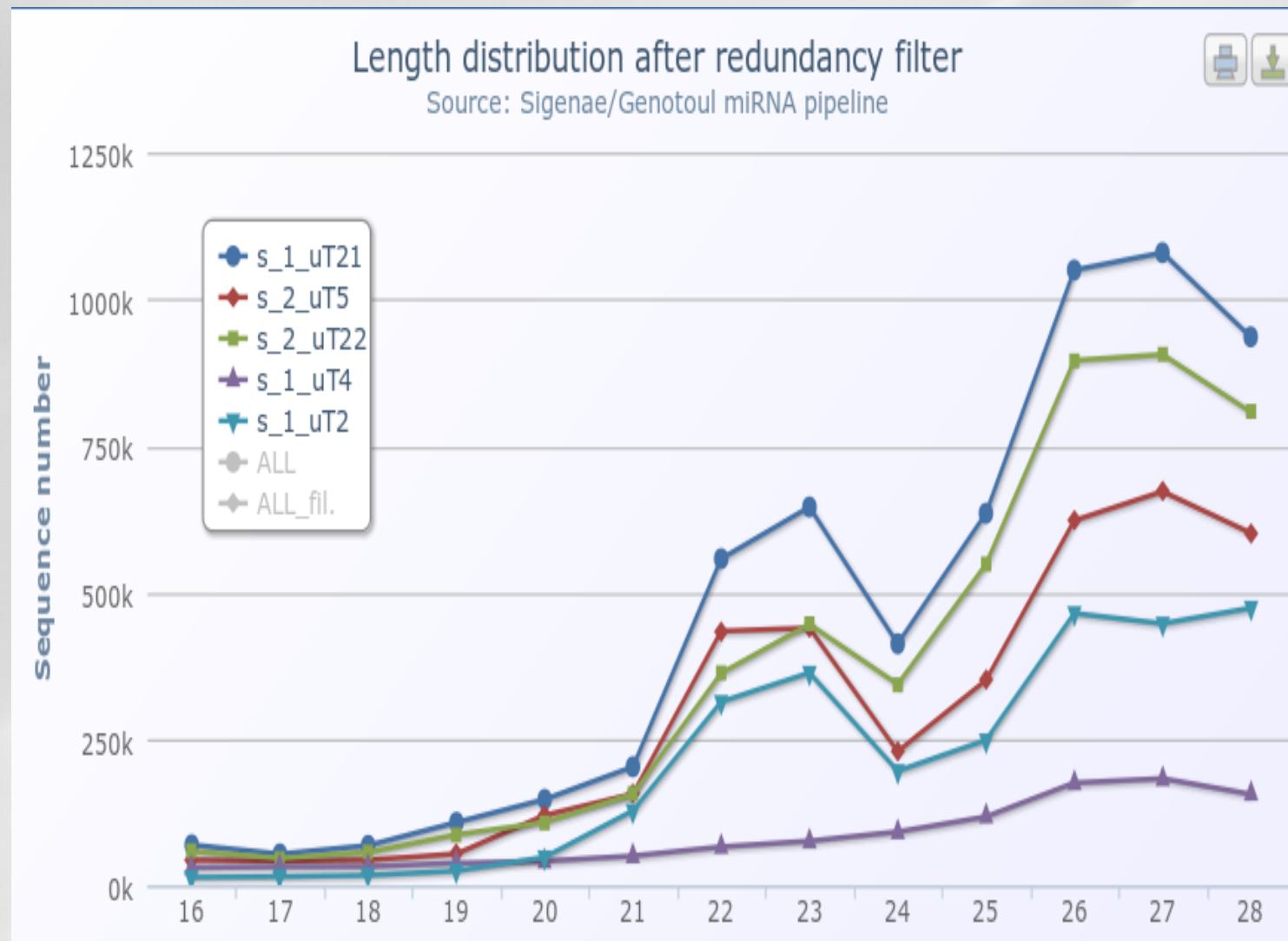
Known miRNA
Quantification matrix

- **Removing identical reads**

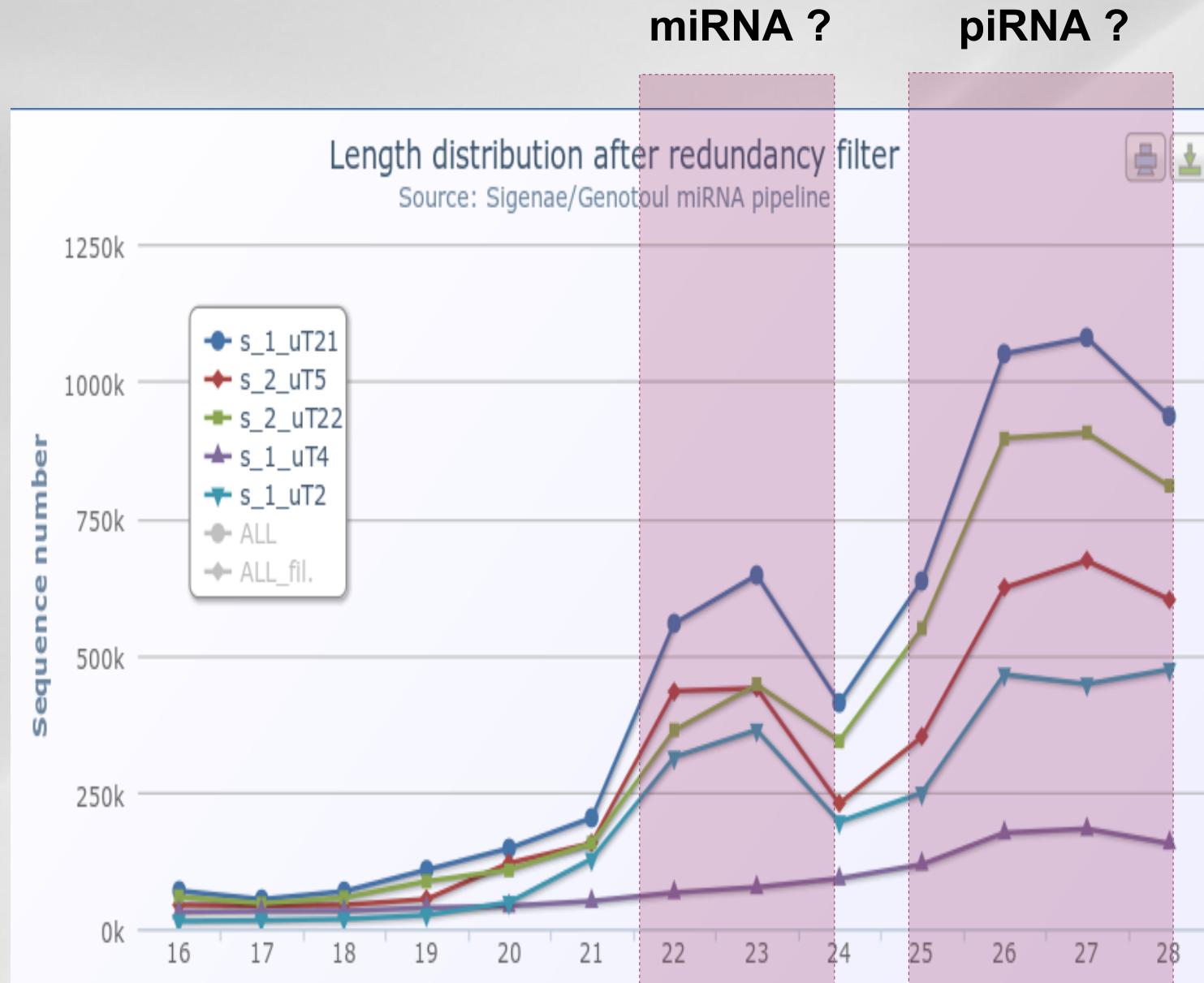
- save computational time
- useless to keep all the read
- Keep the number of occurrence for each read

```
...
AAATGAATGATCTATGGACAGCA          2
AAATGAATGATCTATGGACAGCAG         38
AAATGAATGATCTATGGACAGCAGA        2
AAATGAATGATCTATGGACAGCAGAAAG     1
AAATGAATGATCTATGGACAGCAGCAGC      51
AAATGAATGATCTATGGACAGCAGCA       82
AAATGAATGATCTATGGACAGCAGCAA      5
AAATGAATGATCTATGGACAGCAGCAGAAA    2
AAATGAATGATCTATGGACAGCAGCAGAAC    3
AAATGAATGATCTATGGACAGCAGCAGCAAG   57
AAATGAATGATCTATGGACAGCAGCAGCAG     2
AAATGAATGATCTATGGACAGCCGC          1
AAATGAATGATCTATGGACAGGGCAGCA        1
...
...
```

Remove redundancy



Remove redundancy



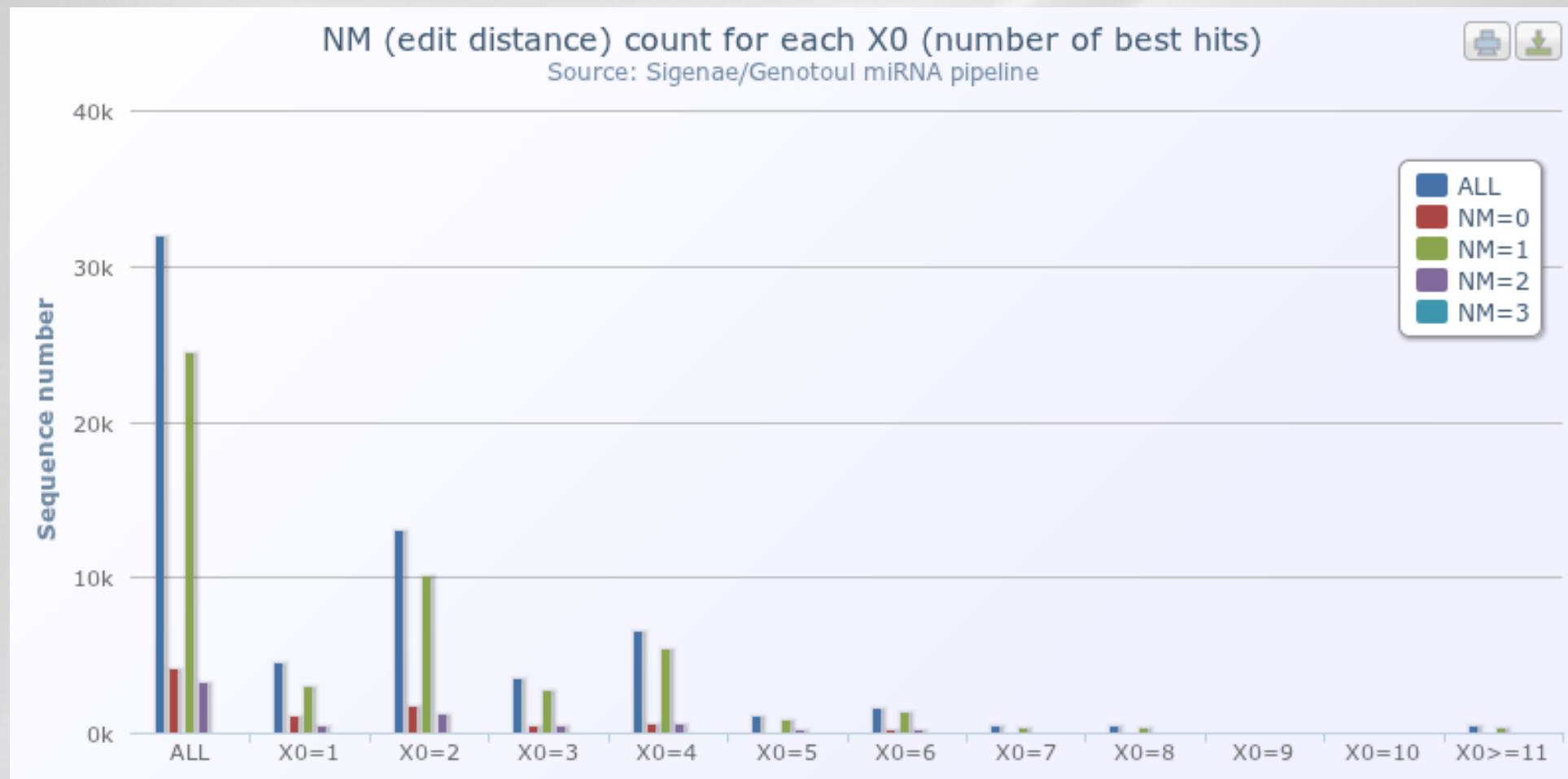
- More differences between piRNAs than with miRNAs ?

Mapping the reads

- Blat <http://genome.ucsc.edu/cgi-bin/hgBlat>
- Blast <http://blast.ncbi.nlm.nih.gov/Blast.cgi>
- Gmap <http://www.gene.com/share/gmap/>
- **Bowtie** <http://bowtie-bio.sourceforge.net/index.shtml>
- **BWA** <http://bio-bwa.sourceforge.net>
- ...

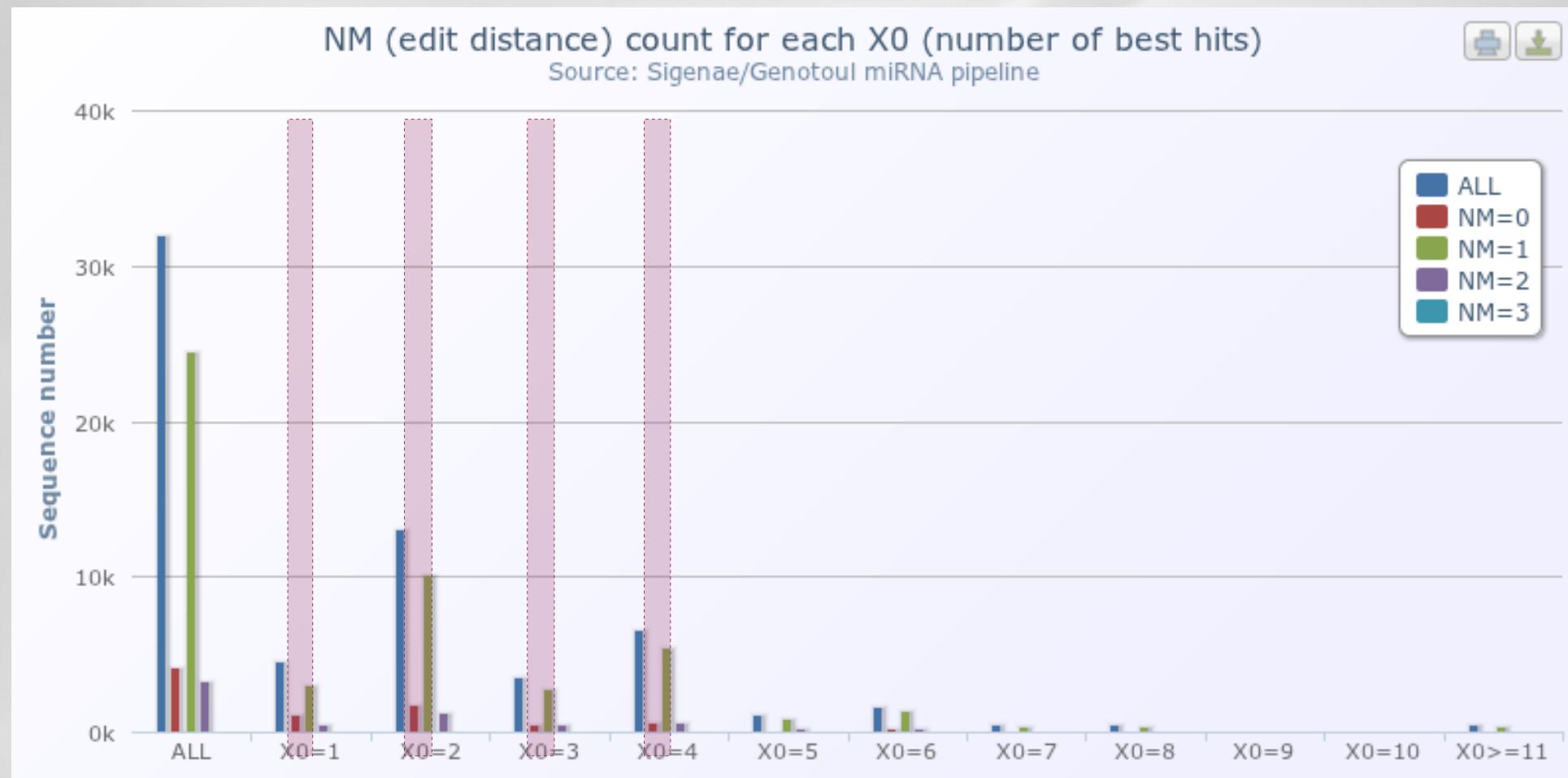
Mapping the reads

- Alignment of annotated reads



Mapping the reads

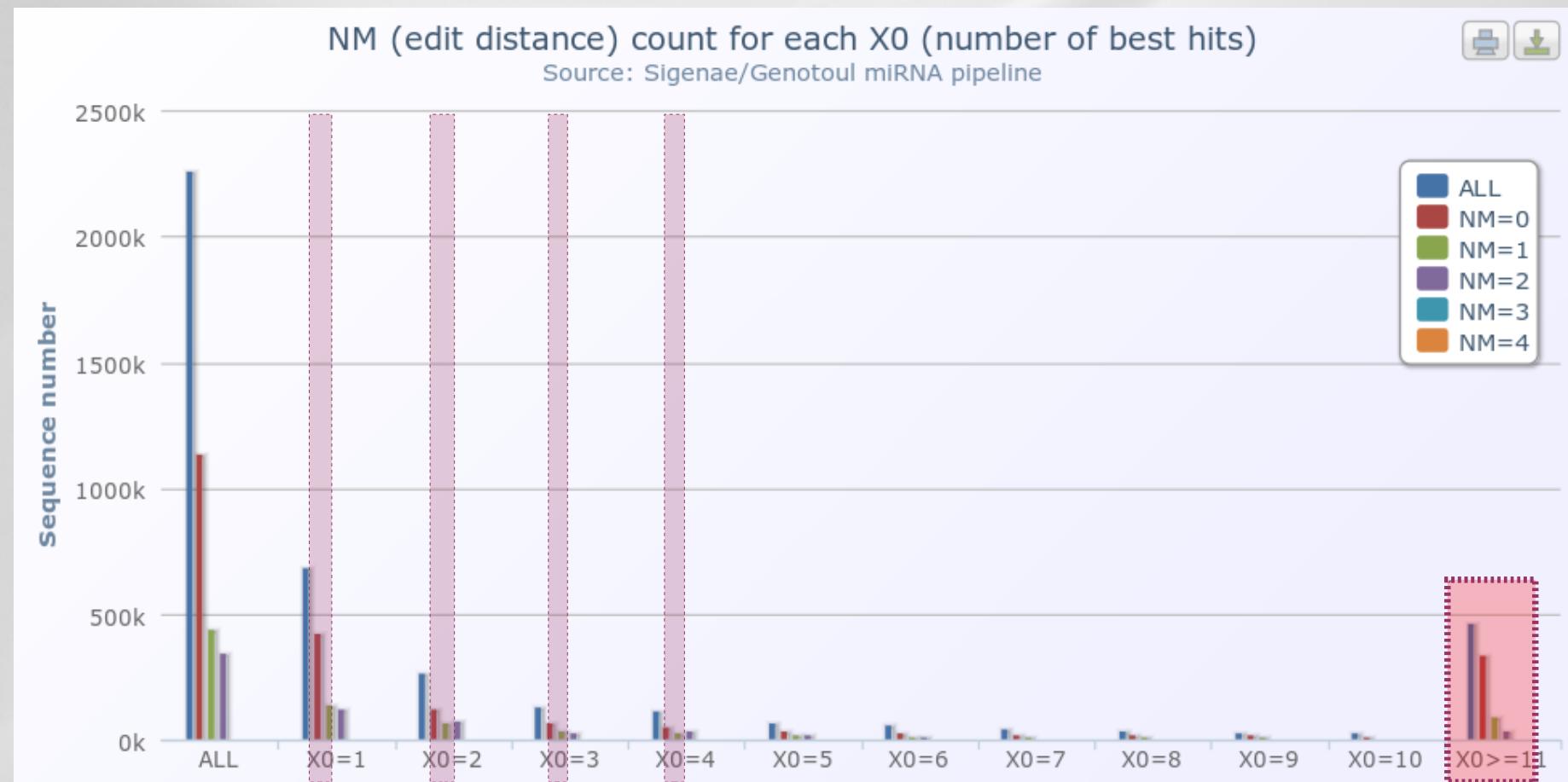
- Alignment of annotated reads



→ keep reads aligned the most at 4 positions with 0 or 1 error

Mapping the reads

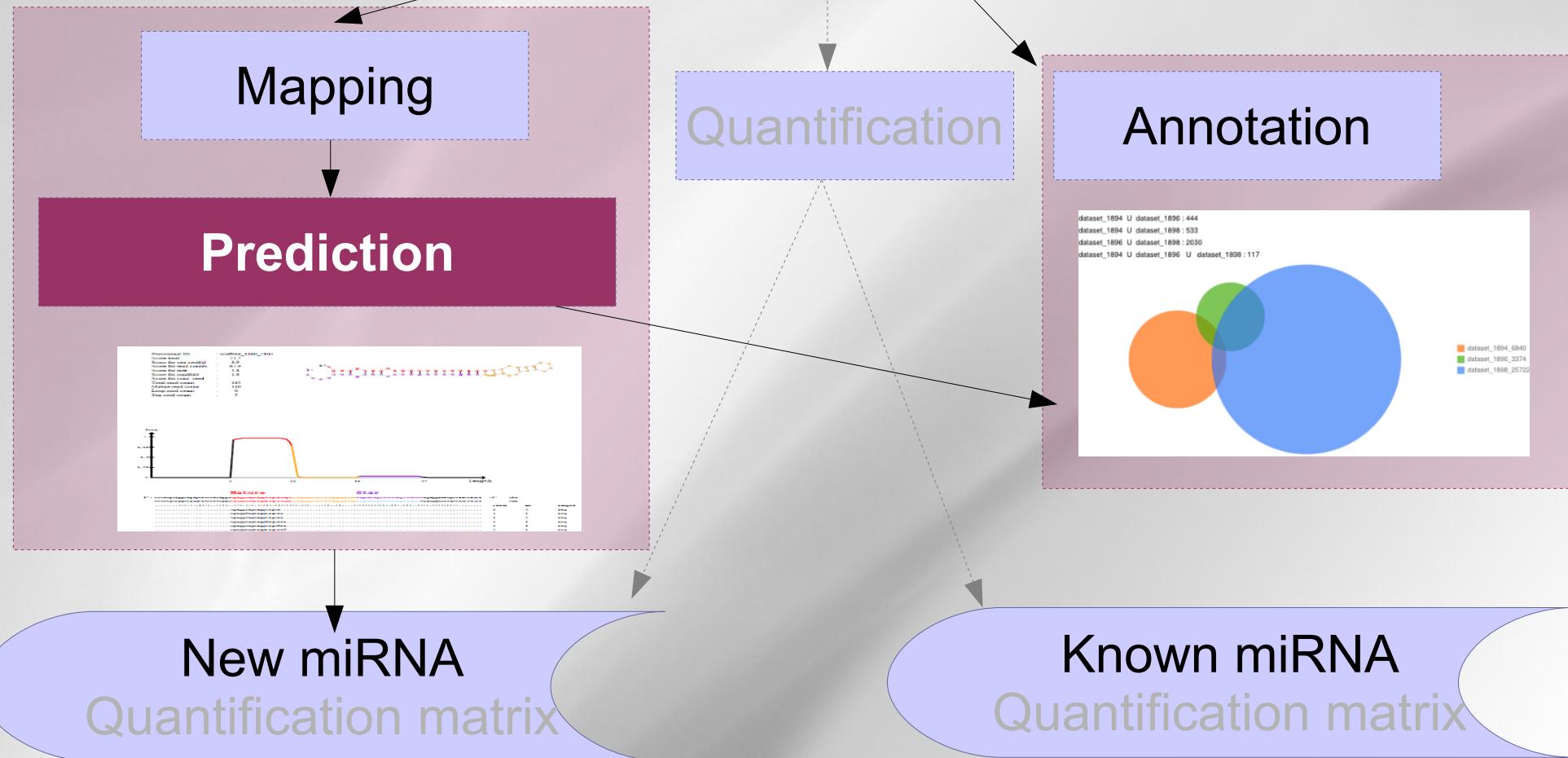
- Alignment of all reads



→ keep reads aligned the most at 4 positions with 0 or 1 error

small RNAseq pipeline

with reference

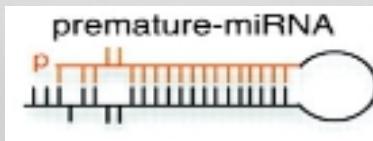


Exercices:

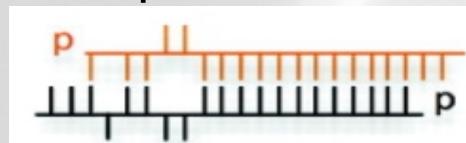
- **Mapping the reads with miRDeep2**
 - **Using Bowtie for mapping**
 - **miRDeep2-core for miRNA identification**

What should we retain for data analysis ?

- Pre-miRNA information:



- Hairpin structure of the pre-miRNA
- Pre-miRNA localisation (coding/non coding TU intronic/exonic)
- Presence of cluster
- Size of the pre-miRNA
- miRNA-5p and miRNA-3p information:



- Existence of both miRNA-5p and miRNA-3p
- Sequence conservation
- Overhang (around 2 nt) related to Drosha and Dicer cuts
- Size of miRNA-5p and miRNA-3p
- Overexpression of one of the miRNA-5p and miRNA-3p
- Existence of other products in sRNAseq data

- Precise excision of a 21-22mer is typical of microRNA
 - less represented reads are products of Dicer errors and sequencing/sample preparation artifacts

GAGACTGGAGTGCAGCCAAGGATGACTTGCCGGATTACATATAGAGTGGAATGA	
<u>CAGCCAAGGATGACTTGCCGG</u>	675
CAGCCAAGGATGACTTGCCGG	26
AGCCAAGGATGACTTGCCGG	9
CAGCCAAGGATGACTTGCCGGAA	8
CAGCCAAGGATGACTTG	2
CAGCCAAGGATGACTTGCCGGA	2
CAGCCAAGGATGACTTGC	1

- Once the reads mapped



- Identify all contiguous read regions



Prediction

- Identify all contiguous read regions



- miRNA precursors have a characteristic secondary structure
 - The detection of a microRNA* sequence, opposing the most frequent read in a stable hairpin (but shifted by 2 bases), is sufficient to diagnose a microRNA.

Mir-30	CTGTAAACATCCTTGACTGGAAGCTGG ***** (((((****((**((((((*****)))))))))))*** 00000000011111111122222222233333333444444444555555555666666666 1234567890123456789012345678901234567890123456789012345678
2	***** CTTTCA
60	***** CTTTCA
8	***TAAACATCCTTGACTGGAAGCTGG*** *****
10	***TAAACATCCTTGACTGGAAGCTG*** *****
89	***TAAACATCCTTGACTGGAAGCT*** *****
297	***GTAAACATCCTTGACTGGAAGCT*** *****
1677	**GTAAACATCCTTGACTGGAAGC*** *****
2	**GTAAACATCCTTGACTGGAAGCTG*** *****
459435	*TGTAACATCCTTGACTGGAAGC* ***** *****
30331	*TGTAACATCCTTGACTGGAAG***** *****
40391	*TGTAACATCCTTGACTGGAAGCT***** *****
17	CTGTAAACATCCTTGACTGGAAGCT***** *****
259	CTGTAAACATCCTTGACTGGAAGC***** *****
21	CTGTAAACATCCTTGACTGGAAG***** *****
2	CTGTAAACATCCTTGACTGGAAG***** 1234567890123456789012345678901234567890123456789012345678 00000000011111111122222222233333333444444445555555566666666666

N	N	N	N	N
N	N	N	N	N
G	G	G	G	G
U	G	C	A	N
A	A	C	G	C
G	G	G	A	A
U	U	U	U	G
C	C	C	G	C
A	A	C	A	C
C	C	C	A	C
A	A	A	A	C
U	U	U	U	G
G	G	G	A	A
U	U	U	A	A
C	C	C	C	U

Prediction

- Extend and fold read regions

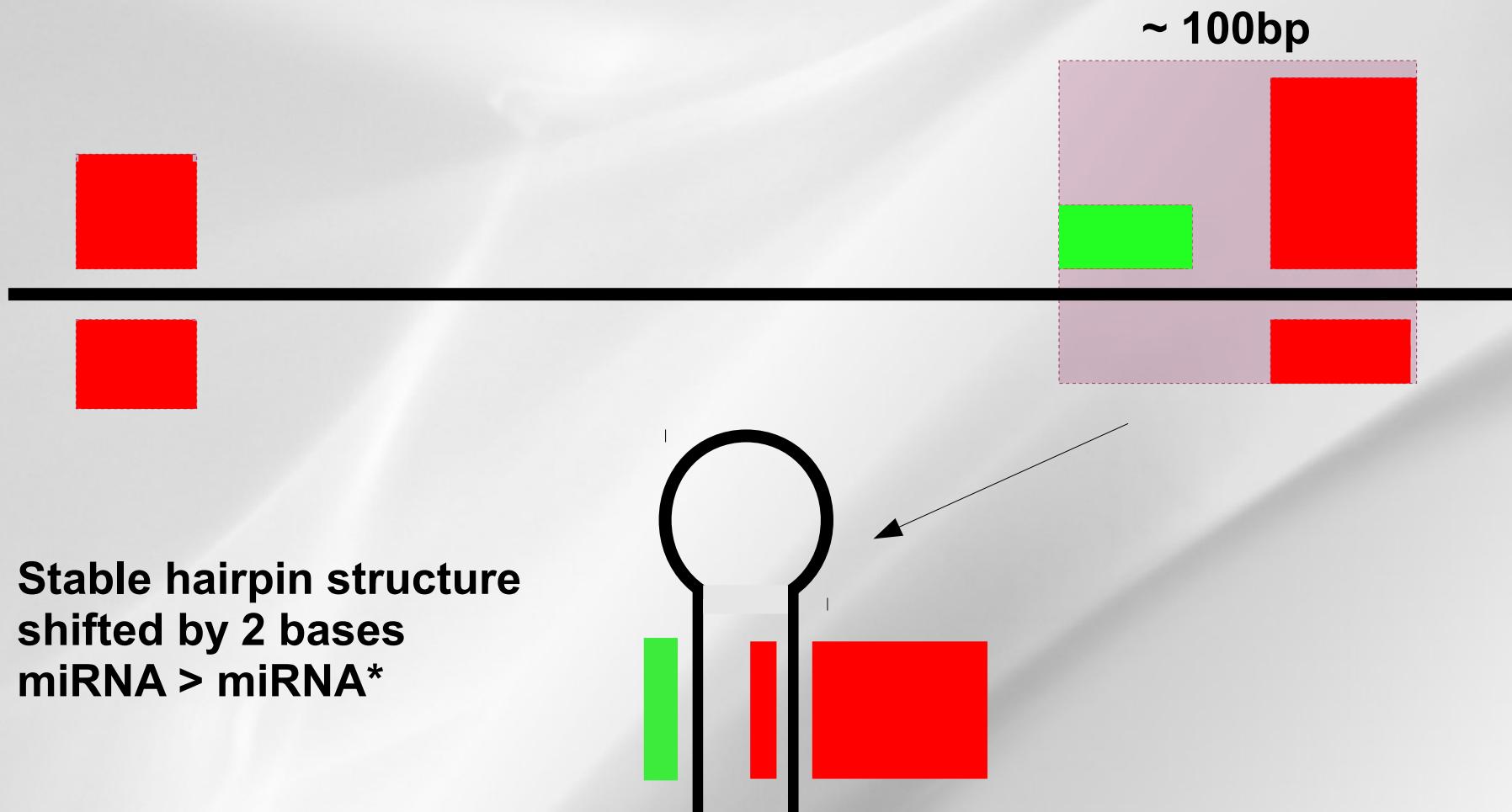


Prediction

- Extend and fold read regions



- Extend and fold read regions



Prediction

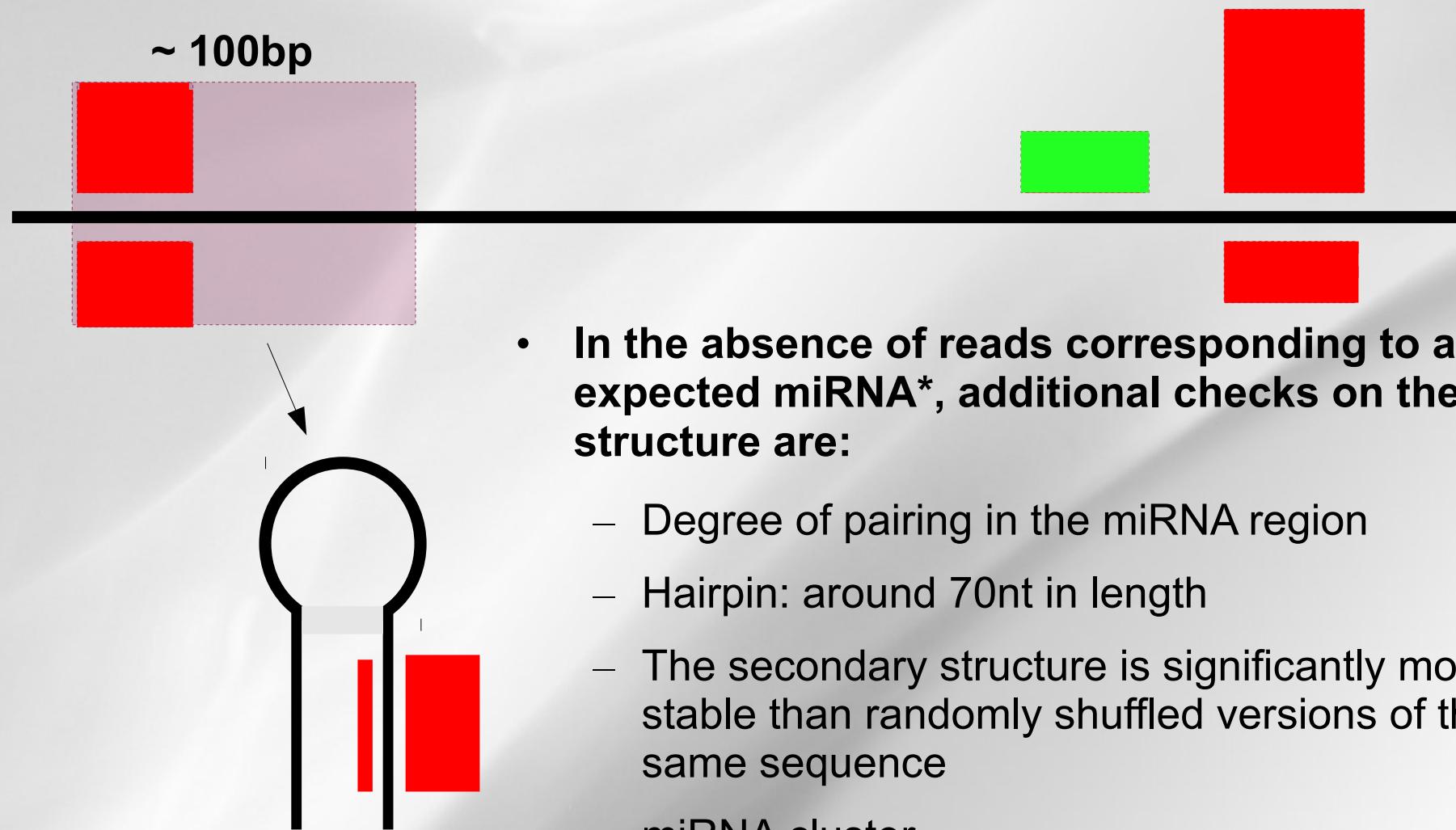
- Extend and fold read regions



OR

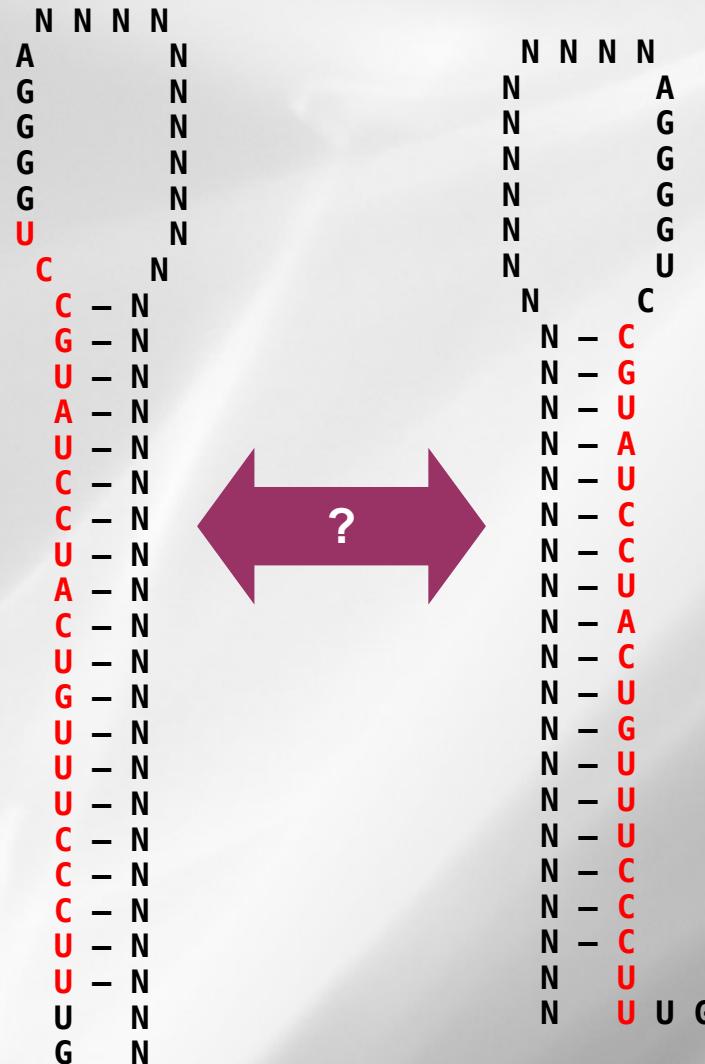


- Extend and fold read regions



- Degree of pairing in the miRNA region
- Hairpin: around 70nt in length
- The secondary structure is significantly more stable than randomly shuffled versions of the same sequence
- miRNA cluster

- Which one should be used ?



Mir-204	GTTCCCTTGTACCTATGCC
2	*****TTTGTACCTATGCCGGAGA
3	***TCCCTTGTACCTATGCCG***
3	**TCCCTTGTACCTATGCCGGAG*
37033	**TCCCTTGTACCTATGCC*****
1597	**TCCCTTGTACCTATGCCG***
2	**TCCCTTGTACCTATGCCGG***
6561	*TTCCCTTGTACCTATGCC*****
611	*TTCCCTTGTACCTATGCC*****
2	GTTCCCTTGTACCTATGCC*****
3	GTTCCCTTGTACCTATGCC*****

W68-W76 *Nucleic Acids Research*, 2009, Vol. 37, Web Server issue
doi:10.1093/nar/gkp347

Published online 11 May 2009

miRanalyzer: a microRNA detection and analysis tool for next-generation sequencing experiments

Michael Hackenberg¹, Martin Sturm², David Langenberger^{3,4},
Juan Manuel Falcón-Pérez⁵ and Ana M. Aransay^{1,*}

¹Functional Genomics Unit, CIC bioGUNE, CIBERehd, Technology Park of Bizkaia, 48160 Derio, Bizkaia, Spain,
²Institute for Bioinformatics and Systems Biology, German Research Center for Environmental Health, Ingolstädter
Landstrasse 1 D-85764 Neuherberg, ³Department of Genome-Oriented Bioinformatics, Wissenschaftszentrum

Published online 16 May 2010

Nucleic Acids Research, 2010, Vol. 38, Web Se

DSAP: deep-sequencing small RNA analysis

Published online 12 September 2011

Nucleic Acids Research, 2012, Vol. 40, No. 1 37–52
doi:10.1093/nar/gkr688

miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades

Marc R. Friedländer¹, Sebastian D. Mackow

BIOINFORMATICS APPLICATIONS NOTE

Sequence analysis

CPSS: a computational platform for the analysis of deep sequencing data

Yuanwei Zhang^{1,†}, Bo Xu^{1,†}, Yifan Yang², Rongjun Ban³, H

Howard J. Cooke^{1,4}, Yu Xue^{5,*} and Qinghua Shi^{1,*}

¹Hefei National Laboratory for Physical Sciences at Microscale and School of and Technology of China, Hefei 230027, China, ²Department of Statistics, Ur 40506, USA, ³Department of Computer Science & Technology, Nanjing University Genetics Unit, IGMM, University of Edinburgh, Edinburgh EH4 2XU, UK, and ⁴Huazhong University of Science and Technology, Wuhan 430074, China

Associate Editor: Ivo Hofacker

<http://www.nature.com/naturebiotechnology>

Discovering microRNAs from deep sequencing data using miRDeep

Marc R Friedländer¹, Wei Chen², Catherine Adamidi¹, Jonas Maaskola¹, Ralf Einspanier³, Signe Knespel¹ & Niklaus Rajewsky¹

The capacity of highly parallel sequencing technologies to detect small RNAs at unprecedented depth suggests their use in systematically identifying microRNAs (miRNAs). However, the identification of miRNAs from the large pool of sequenced transcripts from a single deep sequencing run remains a major challenge. Here, we present an algorithm, miRDeep, which uses a probabilistic model of miRNA

and 454 Life Sciences/Roche, can sequence DNA orders of magnitude faster and at lower cost than Sanger sequencing and are evolving so rapidly that increases in sequencing speed by at least another order of magnitude seem likely over the next few years. Although the Solexa/Illumina system can produce ~32 million sequencing reads in one run, read length is currently limited to 35 bp. In contrast, the current 454 platform yields reads up to 200 bases each, although the number of reads

DOI 10.1007/s11103-012-9885-2

shortran: A pipeline for small RNA-seq data analysis

Vikas Gupta^{1,2}, Katharina Markmann¹, Christian N. S. Pedersen², Jens Stougaard¹ and Anders Andersen^{1,*}

¹Centre for Carbohydrate Recognition and Signalling, Department of Molecular Biology and Genetics, Aarhus Gustav Wieds Vej 10, 8000 Aarhus C, Denmark and ²Bioinformatics Research Centre, Aarhus University, C 8, 8000 Aarhus C, Denmark

Prediction Existing software

BMC Bioinformatics



Open Access

miRExpress: Analyzing high-throughput sequencing data for profiling microRNA expression

Wei-Chi Wang¹, Feng-Mao Lin¹, Wen-Chi Chang^{1,5}, Kuan-Yu Lin^{2,3}, Hsien-Da Huang^{*1,4} and Na-Sheng Lin^{*2,3}

Address: ¹Institute of Bioinformatics of Biotechnology, National Chen Sinica, Nankang, Taipei 11529, T Hsin-Chu 300, Taiwan, Republic of China
Email: Wei-Chi Wang - cancer@bi.sinica.edu.tw

Hendrix et al. *Genome Biology* 2010, 11:R39
<http://genomebiology.com/2010/11/4/R39>



METHOD

Open Access

miRTRAP, a computational method for the systematic identification of miRNAs from high throughput sequencing data

NATURE BIOTECHNOLOGY VOLUME 26 NUMBER 4 APRIL 2008

NOTE

Vol. 26 no. 20 2010, pages 2615–2616
doi:10.1093/bioinformatics/btq493

Advance Access publication August 27, 2010

ep sequencing analysis

ov², Gideon Dror², Eran Halperin^{3,4}

edicine, Tel Aviv University, ²The Academic ice Institute, Berkeley, CA, USA and ⁴School Biotechnology, George Wise Faculty of Life

miRDeepFinder: a miRNA analysis tool for deep sequencing of plant small RNAs

Fuliang Xie · Peng Xiao · Dongliang Chen · Lei Xu · Baohong Zhang

Prediction

Existing software

- Basic features
 - Availability (web/executable)
 - Computing resources (time, memory)
 - Reads pre-processing
 - Mapping
 - Identification

Briefings in Bioinformatics Advance Access published March 24, 2012
 BRIEFINGS IN BIOINFORMATICS, page 1 of 10
 doi:10.1093/bib/bbs010

Detecting miRNAs in deep-sequencing data: a software performance comparison and evaluation

Vernell Williamson, Albert Kim, Bin Xie, G. Omari McMichael, Yuan Gao and Vladimir Vladimirov
 Submitted: 9th December 2011; Received (in revised form): 21st February 2012

Table 2: Basic features of popular software used to predict miRNA from deep-sequencing data

Accessible	Read pre-processing	Target genomes	Mapping algorithm	Functions	Predictions based on	Location	Program
Executable requires in-house computational resources.	Provides script that eliminates redundancy. Tag removal/processing must be done by user prior to analysis.	Flexible, Human (GRCh37).	Flexible, Oligomap (v1) Bowtie (v2).	Novel, known miRNA prediction. Status of predictions (novel/known) must be determined by the user.	Bayesian probability, focus on traditional steps of biogenesis.	http://www.mdc-berlin.de/en/research/research-teams/	MiRDeep/miRDeep2
Web based	Accepts two multifasta format and file with read and counts. Tag must be removed by user.	Seven genomes (human, fruit fly, rat, mouse, dog, nematode, and zebra fish), fixed choice over version.	Fixed, BowTie. User can set the number of acceptable mismatches (<2).	Novel, Known miRNA prediction.	Posterior probability (threshold > 0.95). Reads are mapped against target genome, miRBase, and other non-coding databases.	http://web.bioinformatics.cicbiogune.es/microRNA/miRAnalyser.php	MirAnalyzer
Web-based	Accepts read/counts format like miRAnalyzer. Adapter sequences can be left intact	Multiple genomes, fixed choice over version	Fixed, cluster approach, Uses SuperMatcher to increase speed	Known miRNA prediction, species distribution, expression level	Degree to which reads match known examples. Known miRNAs are compared to miRBase	http://dsap.cgu.edu.tw/	DSAP

Prediction

Existing software

- Reads pre-processing
 - Adaptators trimming
 - Redundancy
 - Repeats
 - Other ncRNA
 - Size of the mature miRNA (min/max)

Prediction

Existing software

- Mapping
 - Size and region of the read
 - Number of locations
 - Considered
 - Reported
 - Error(s) consideration in mapping
 - Quality of the read

Prediction

Existing software

- Precursor identification
 - Length and bounds of the theoretical sequence (folding)
 - Alignment of the read against known miRNA
- Post procesing step: assessment of the potential miRNA
 - Differents methods: SVM, bayesian statistics based score, combinatorial rules...
 - Location of the read on the precursor
 - 2 nt overhang of the mature miRNA/precursor
 - Accuracy of the folding (HP structure, energy, Z-score...)

Existing software miRDeep & miRDeep2

NATURE BIOTECHNOLOGY VOLUME 26 NUMBER 4 APRIL 2008

Discovering microRNAs from deep sequencing data using miRDeep

Marc R Friedländer¹, Wei Chen², Catherine Adamidi¹, Jonas Maaskola¹, Ralf Einspanier³, Signe Knespel¹ & Nikolaus Rajewsky¹

<http://www.nature.com/naturebiotechnology>

The capacity of highly parallel sequencing technologies to detect small RNAs at unprecedented depth suggests their value in systematically identifying microRNAs (miRNAs). However, the identification of miRNAs from the large pool of sequenced transcripts from a single deep sequencing run remains a major challenge. Here, we present an algorithm, miRDeep, which uses a probabilistic model of miRNA

and 454 Life Sciences/Roche, can sequence DNA orders of magnitude faster and at lower cost than Sanger sequencing and are evolving so rapidly that increases in sequencing speed by at least another order of magnitude seem likely over the next few years. Although the Solexa/Illumina system can produce ~32 million sequencing reads in one run, read length is currently limited to 35 bp. In contrast, the current 454 platform yields reads up to 200 bases each, although the number of reads

Published online 12 September 2011

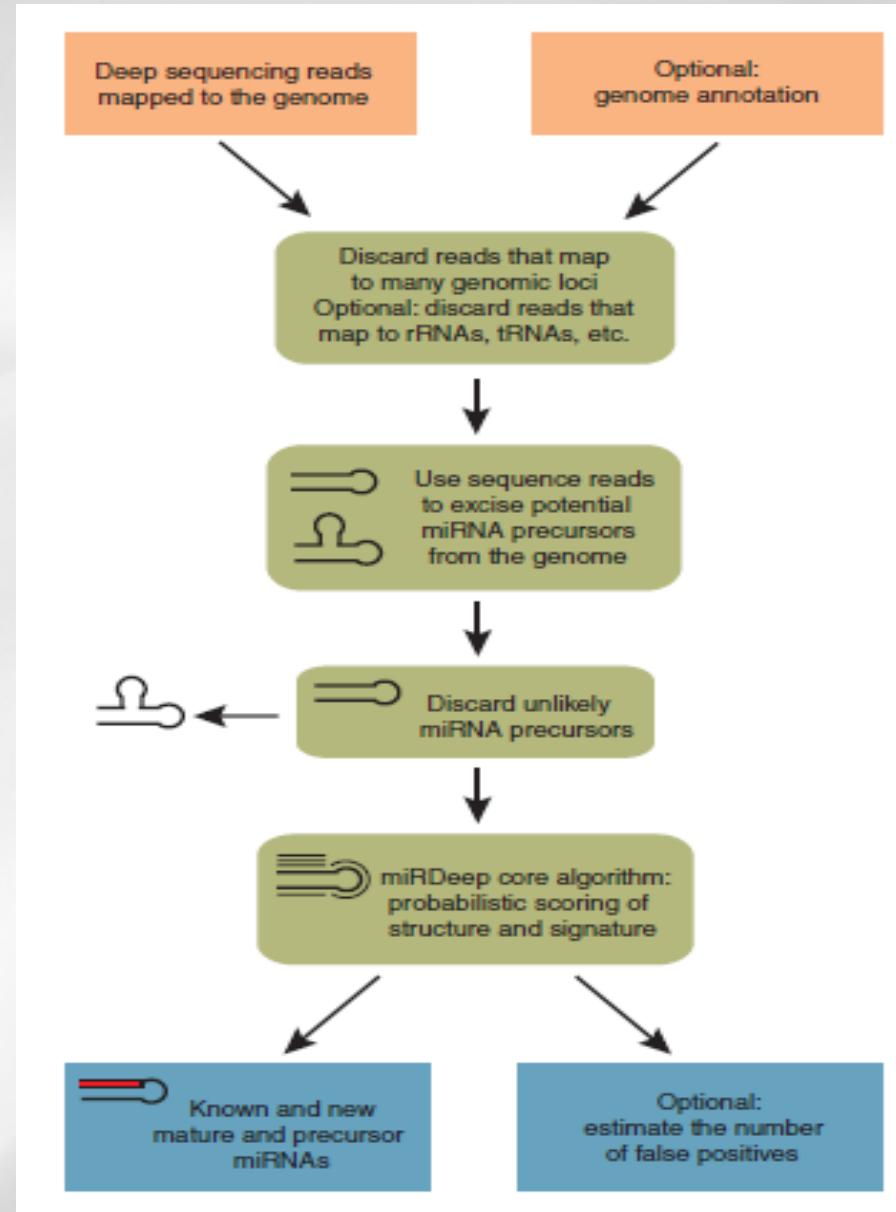
Nucleic Acids Research, 2012, Vol. 40, No. 1 37–52
doi:10.1093/nar/gkr688

miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades

Marc R. Friedländer¹, Sebastian D. Mackowiak¹, Na Li², Wei Chen² and Nikolaus Rajewsky^{1,*}

¹Laboratory for Systems Biology of Gene Regulatory Elements and ²Laboratory for New Sequencing Technology, Berlin Institute for Medical Systems Biology at the Max-Delbrück-Center for Molecular Medicine, Berlin-Buch 13125, Germany

Existing software miRDeep2



Existing software

miRDeep

- <Http://www.mdc-berlin.de/rajewsky/miRDeep>
- Seven Perl scripts
- Required dependancies
- Vienna package (Hofacker, NAR, 2003)
- Randfold program (Bonnet et al., Bioinformatics, 2004)
- Megablast (Altschul et al. J. Mol. Biol. 1990)

Existing software

miRDeep

- Mapping
 - Megablast
 - Mismatches tolerated in the last three nucleotides
- Theretical precursor
 - Reads aligned in a 30nt distance
 - Clustering plus 25nt flanks
 - Sequences of more than 140nt discarded
 - Reads not aligned in a 30nt distance
 - Two potential precursors of length 110nt

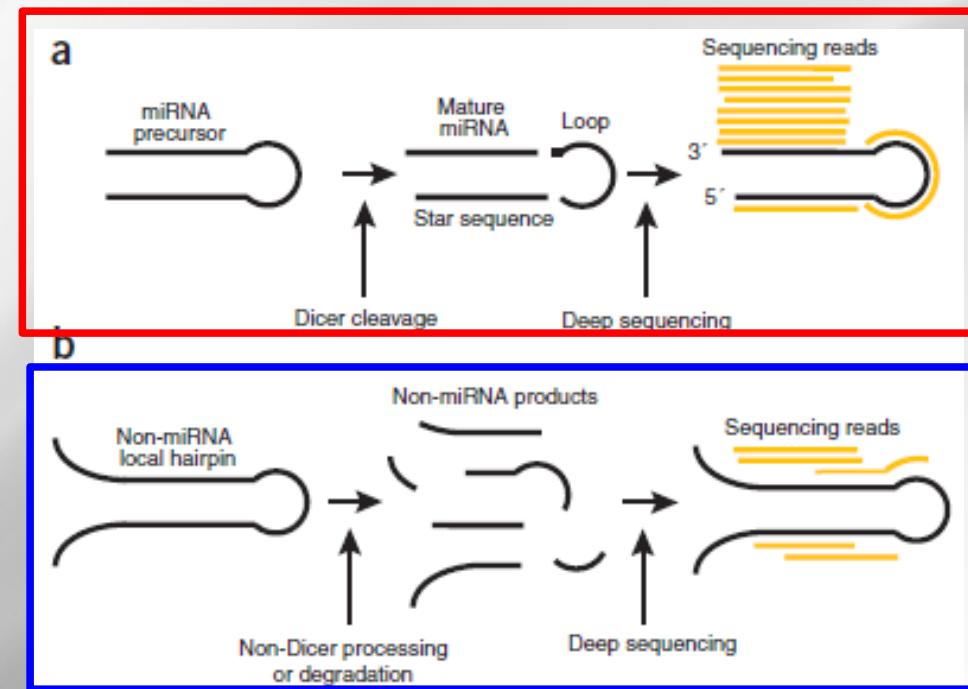
Existing software

miRDeep & miRDeep2

- At least 11 million hairpins in the human genome (Bentwich, FEBS Lett., 2005)
- Most reads originate from loci that are not miRNA genes

Score a characteristic signature of
miRNA compared to **other short RNA products**

- Energetic stability: relative and absolute
- Positions on the precursor and 2nt 3' duplex overhang
- 5' end of the potential mature sequenced conserved in known mature miRNA (option)
- frequencies (over-representation of reads in the mature miRNA)
- Bayesian statistics to score the fit of sequenced RNA to the biological model:
 - $\text{score} = \log(P(\text{pre}/\text{data})/P(\text{bgr}/\text{data}))$



Existing software

miRDeep

- The PERL scripts
 - **blastoutparse.pl**: parse standard blast output into a custom tabular separated format
 - **blastparseselect.pl**: cleans the output of blastoutparse.pl
 - **filter_alignments.pl**: filters the alignments of deep-sequencing reads to a genome
 - **excise_candidates.pl**: cuts out potential precursors
 - **mirdeep.pl**: core algorithm which provides all information on candidate precursors
 - **permute_structure.pl**: do the permutation controls

Existing software miRDeep2

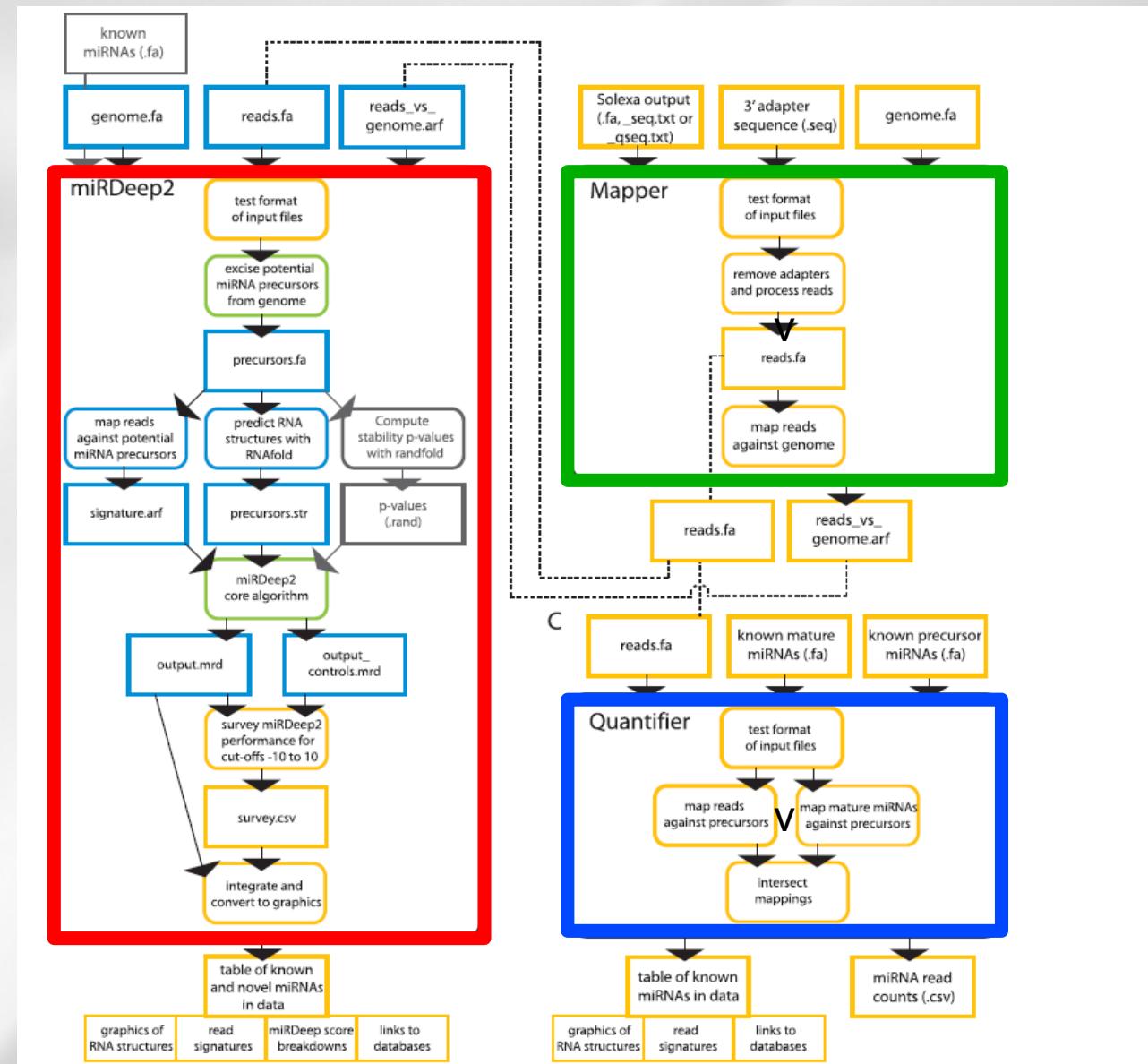
- Improve miRDeep in several ways
 - Easier to launch
 - Faster (all mapping done with Bowtie)
 - Several samples considered simultaneously
 - Increased robustness to identify non canonical miRNA (moRs)
 - Sens/anti-sens reads analyzed separately
 - Graphic output

Existing software

miRDeep2

Three modules

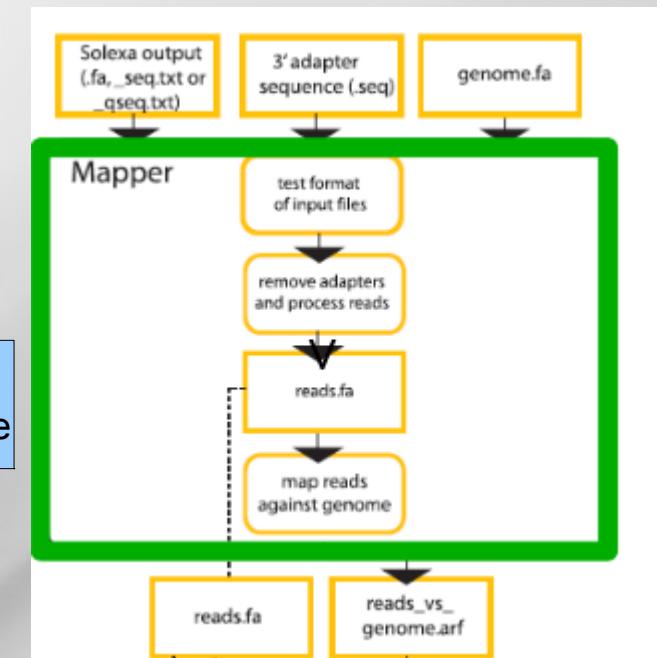
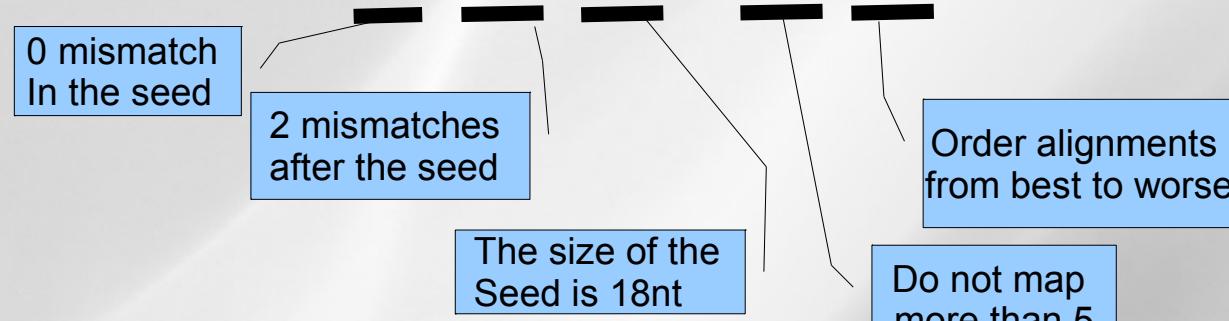
- MiRDeep2
- Mapper
- Quantifier



Existing software miRDeep & miRDeep2

The Mapper module

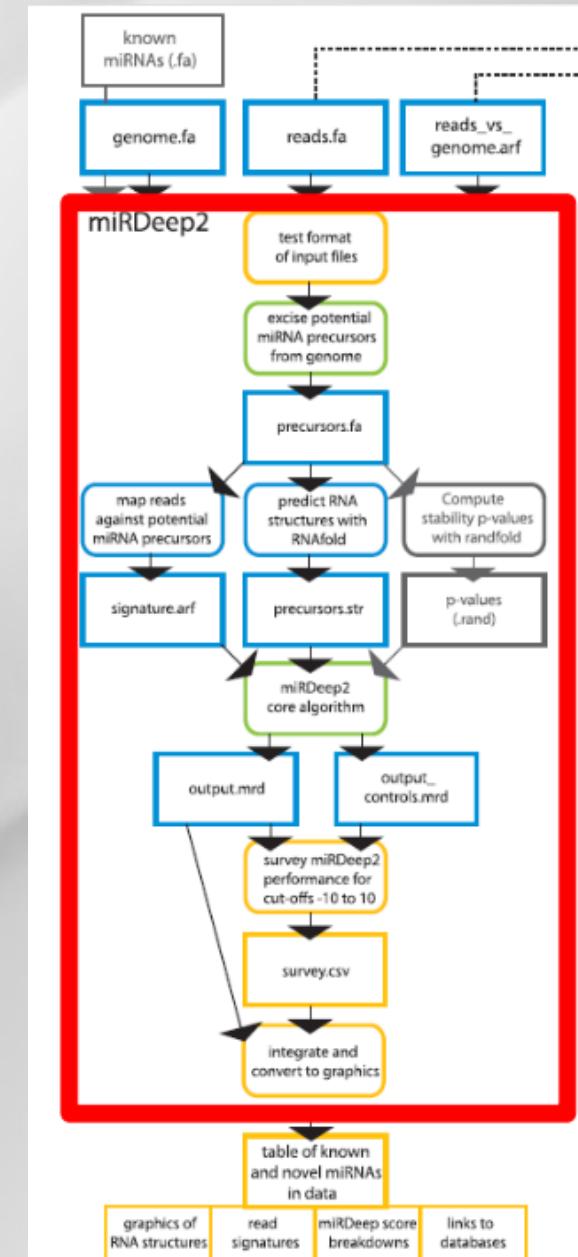
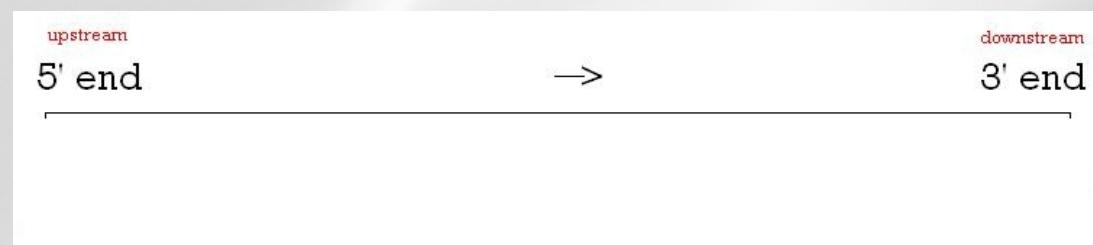
- Reads processing
 - Remove redundancy and keep # occurrences
- Map reads with bowtie (or BWA)
 - Bowtie -f -n 0 -e 80 -l 18 -a -m 5 -best -strata



Existing software miRDeep & miRDeep2

The miRDeep2 module

- Scan of both strands from 5' to 3'

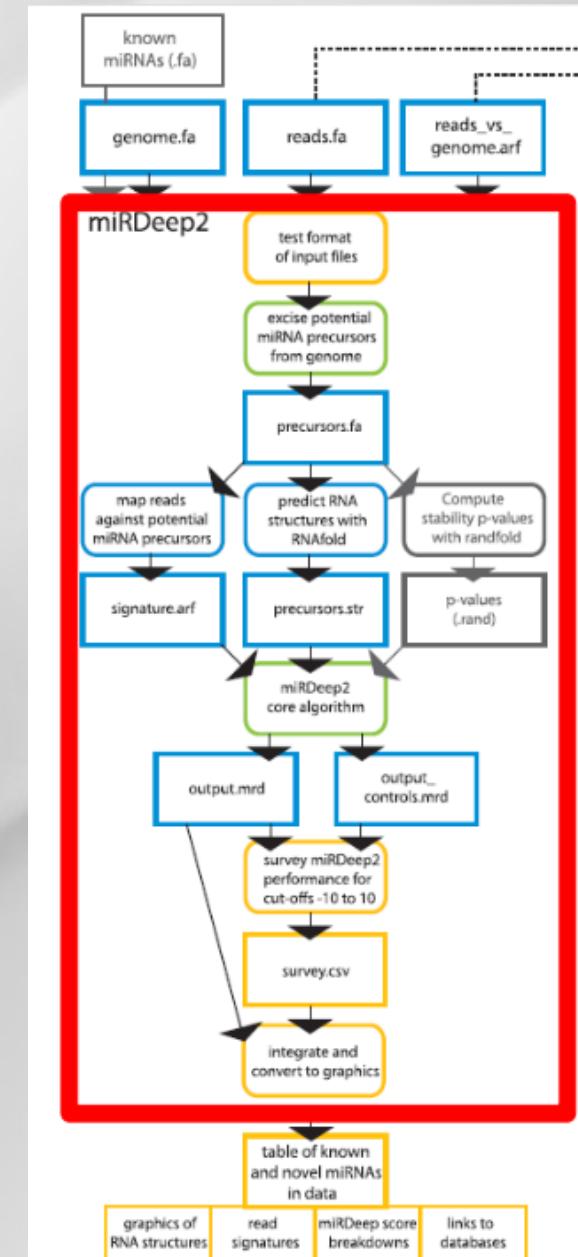
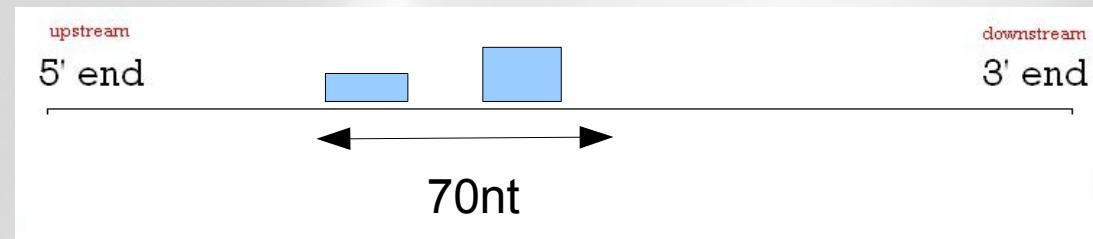


Existing software

miRDeep & miRDeep2

The miRDeep2 module

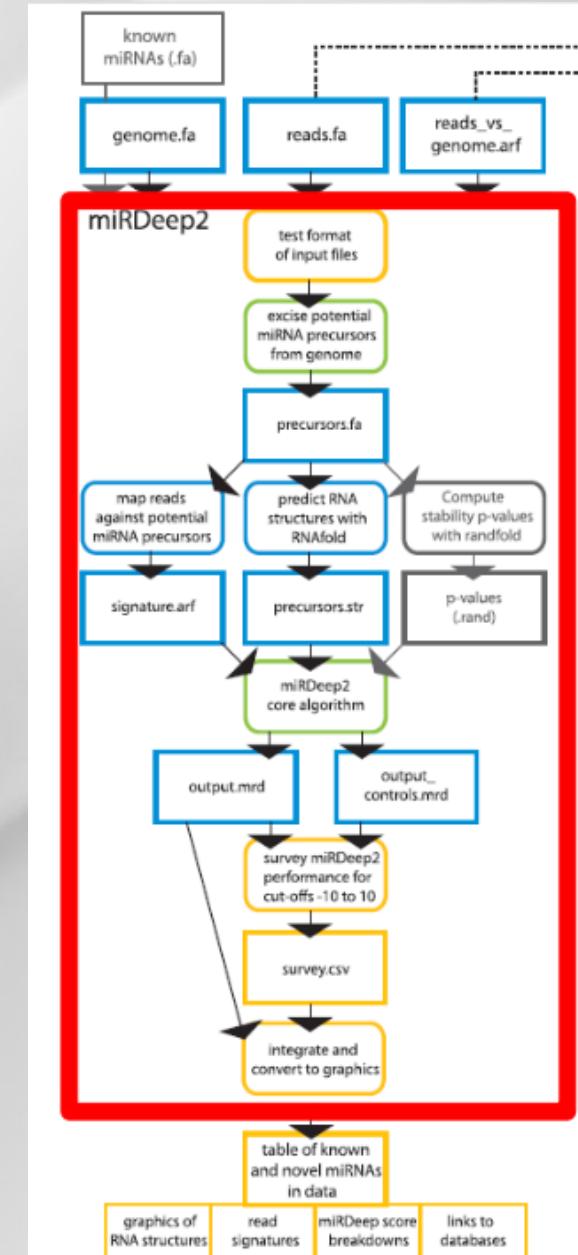
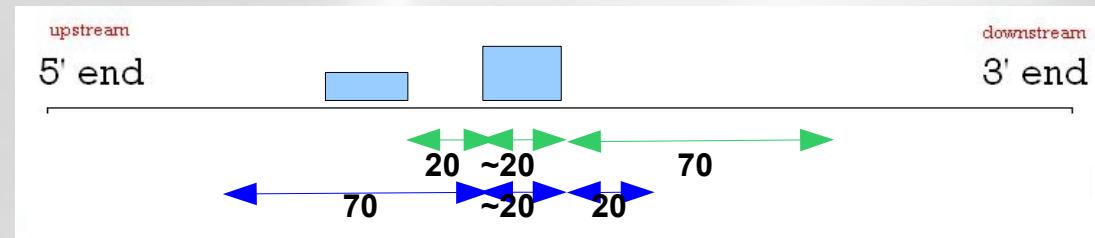
- Scan of both strands from 5' to 3'
 - Search the best stack of reads (height 1 or more) in a distance of 70nt



Existing software miRDeep & miRDeep2

The miRDeep2 module

- Scan of both strands from 5' to 3'
 - Search the best stack of reads (height 1 or more) in a distance of 70nt
 - Excise potential precursors on both sides

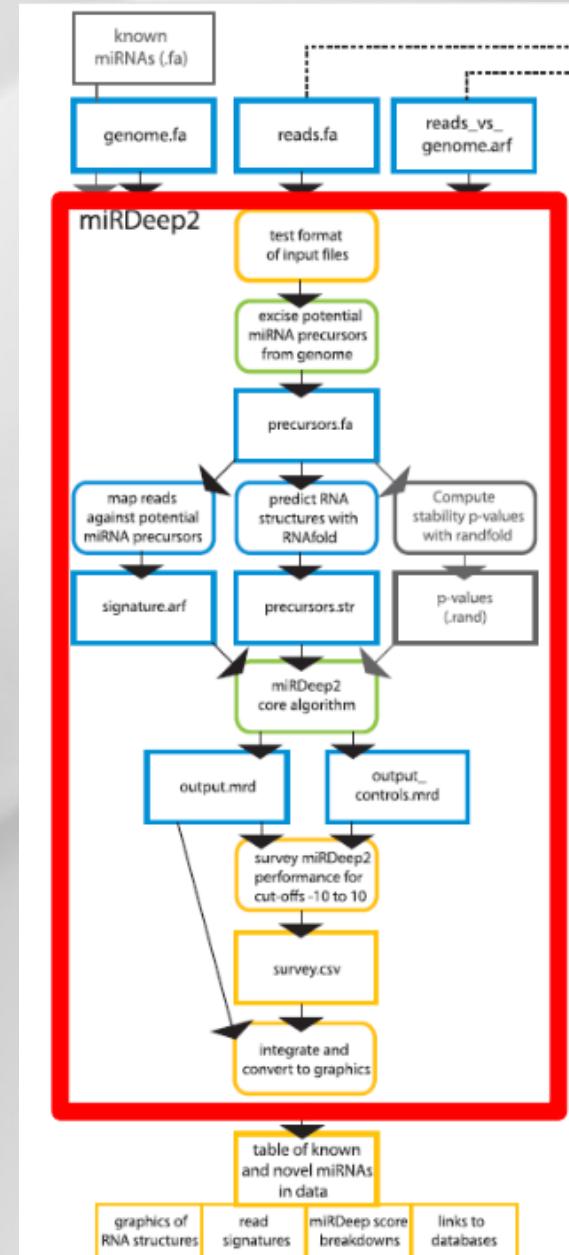


Existing software

miRDeep2

The miRDeep2 module

- Scan of both strands from 5' to 3'
 - Search the best stack of reads (height 1 or more) in a distance of 70nt
 - Excise potential precursors on both sides
 - Go on from 1 nt after the last position excised
- If the number of candidate precursor > 50.000, repeat the process (height of stack = height of stack + 1)
- Prepare the file of precursor signature
 - Align reads against precursors (1 MM allowed)
 - Align known miRNA against precursors (0 MM allowed)
- Evaluation of candidate precursors
 - Fold candidate precursors (RNAlign + Randfold)
 - Unbifurcated hairpins
 - Score the candidates
 - Valid alignment of reads on the precursor
 - 60% of nt in the mature part paired
- Estimation of true positives



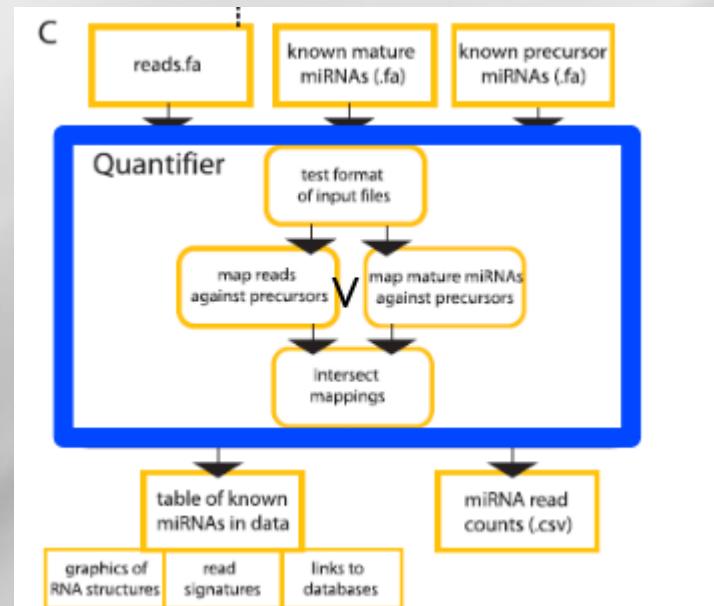
Existing software miRDeep2

The Quantifier module

Identifies and quantifies known mature miRNA given

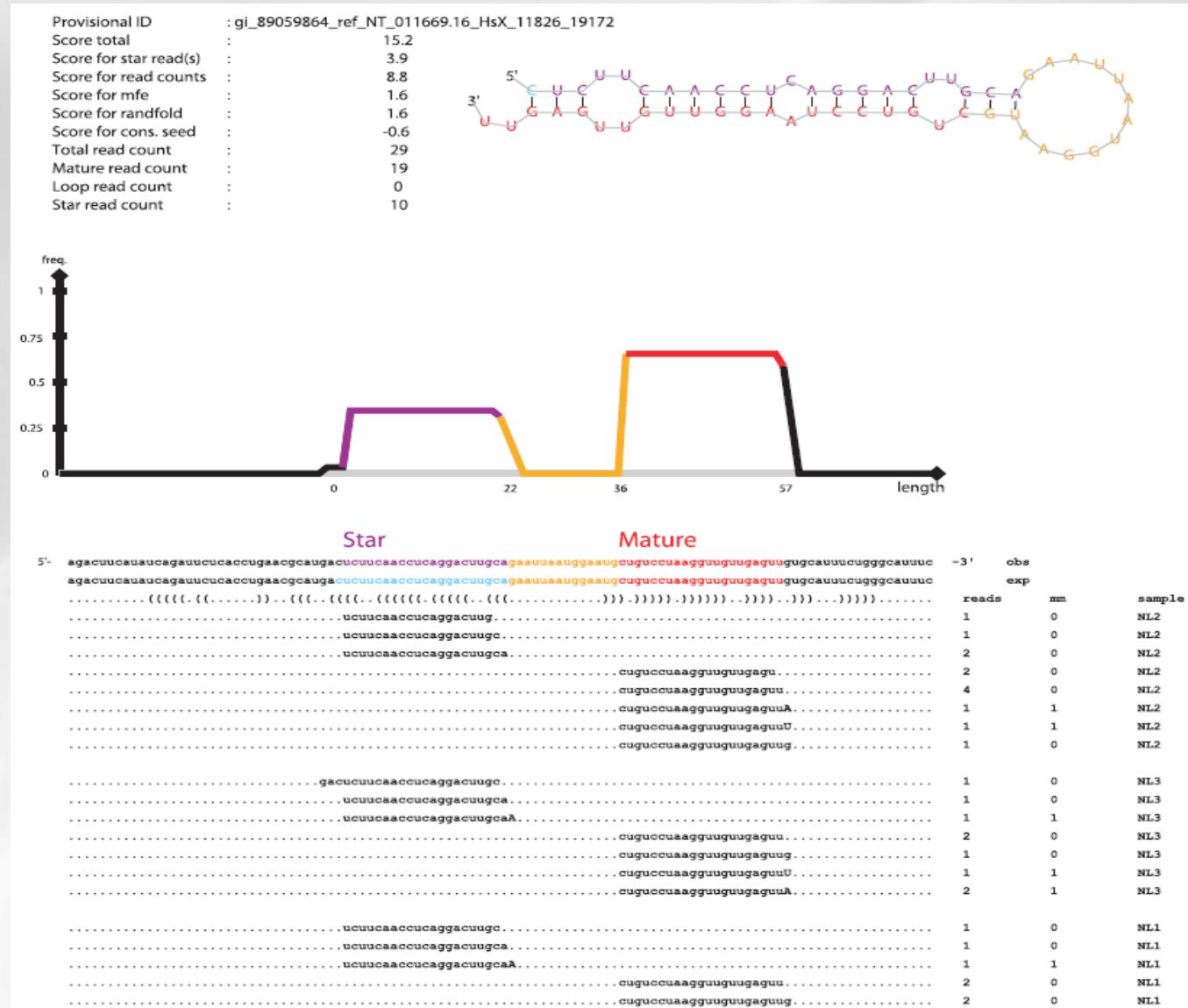
- Know mature miRNA
- Know miRNA precursors

Use Bowtie for miRNA/reads alignment



Existing software

miRDeep2

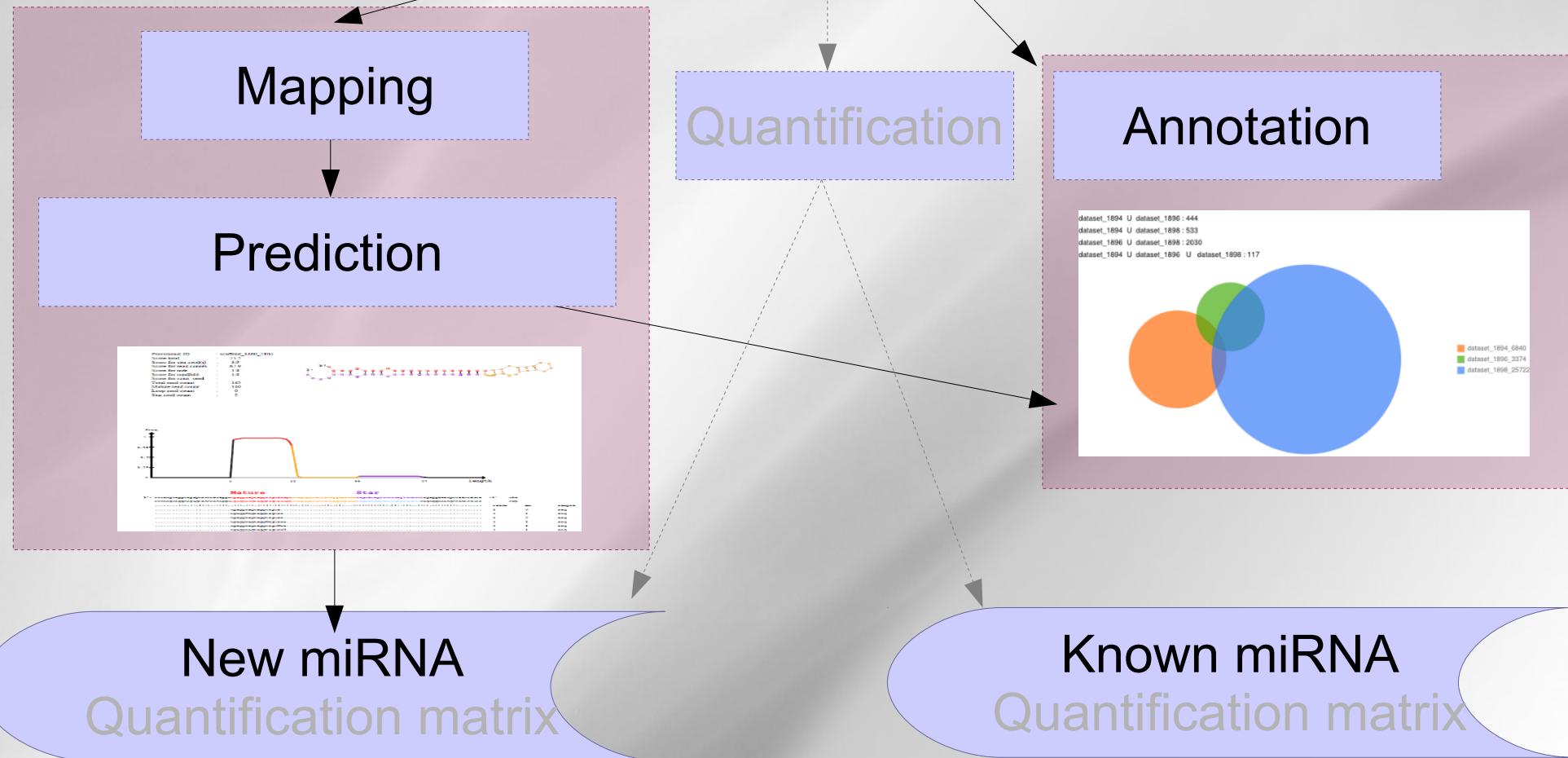


Exercice:

- Back to miRdeep2-core results

small RNAseq pipeline

with reference



- Useful databases:

- miRbase (<http://microrna.sanger.ac.uk/>)



- miRBase::Registry provides names to novel miRNA genes prior to their publication.
- **miRBase::Sequences provides miRNA sequence data, annotation, references and links to other resources for all published miRNAs.**
- miRBase::Targets provides an automated pipeline for the prediction of targets for all published animal miRNAs.

Stem-loop sequence MI0000082

Accession MI0000082
ID hsa-mir-25
Symbol HGNC_MIRN25
Description Homo sapiens miR-25 stem loop

Stem-loop 

[Get sequence](#)

Genome context [? 98915834-993393171](#)

ENST00000349202; intern S.
ENST00000410402; intern 21; NP_06
ENST00000318822; intern 36; MCM
ENST00000362828; intern 54; NP_06

Database links EMBL: A1421746
HGNC: 3169

Related entries mmu-mir-25, msu-mir-25, dm3-mir-25, gpc-mir-25, rno-mir-25
[Show details](#)

Mature sequence MIMAT0000081

Accession MIMAT0000081
ID hsa-miR-25
Sequence 5'- UUUCAGGGUGGUAU - 3'
Get sequence

Evidence experimental; cloned [1-2], Northern [1]

Predicted targets MRBASE: hsa-mir-25
PICTAR-VERT: hsa-mir-25
TARGETSCAN: miR-2532/3267

References 1. Identification of small gene coding for small known miRNAs.
Lagos-Quintana M, Rauhut R, Lendeckel W, Tuschl T
Science 294:853-858(2001).
2. "Altered expression profiles of microRNAs during TPA-induced differentiation of HL-60 cells."
Kasaihama K, Nakamura Y, Koza T
Biochem Biophys Res Commun 322:403-410(2004).

D152–D157 Nucleic Acids Research, 2011, Vol. 39, Database issue
doi: 10.1093/nar/gkq1027

Published online 30 October 2010

miRBase: integrating microRNA annotation and deep-sequencing data

Ana Kozomara and Sam Griffiths-Jones*

Faculty of Life Sciences, University of Manchester, Michael Smith Building, Oxford Road, Manchester, M13 9PT, UK

- Useful databases:
 - miRbase (<http://microrna.sanger.ac.uk/>) 
 - Rfam (<http://rfam.sanger.ac.uk/>)
 - A collection of RNA families
 - Rfam 10.1, June 2011, 1973 families
 - A track now included in the UCSC genome browser
 - Be careful: also contains (not all) miRNA families

D136–D140 *Nucleic Acids Research*, 2009, Vol. 37, Database issue
doi:10.1093/nar/gkn766

Published online 25 October 2008

Rfam: updates to the RNA families database

Paul P. Gardner^{1,*}, Jennifer Daub¹, John G. Tate¹, Eric P. Nawrocki²,
Diana L. Kolbe², Stinus Lindgreen³, Adam C. Wilkinson¹, Robert D. Finn¹,
Sam Griffiths-Jones⁴, Sean R. Eddy² and Alex Bateman¹

¹Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SA, UK, ²Howard Hughes Medical Institute, Janelia Farm Research Campus, Ashburn, Virginia, USA, ³Center for Bioinformatics, Department of Biology, University of Copenhagen, Ole Maaløes Vej 5, DK-2200 Copenhagen N, Denmark and ⁴Faculty of Life Sciences, The University of Manchester, Manchester M13 9PL, UK

- Useful databases:

- miRbase (<http://microrna.sanger.ac.uk/>)
- Rfam (<http://rfam.sanger.ac.uk/>)



- Silva (<http://www.arb-silva.de/>)



- A comprehensive on-line resource for quality checked and aligned ribosomal RNA sequence data.

- SSU (16S rRNA, 18S rRNA)
- LSU (23S rRNA, 28S rRNA)

7188–7196 *Nucleic Acids Research*, 2007, Vol. 35, No. 21
doi:10.1093/nar/gkm864

Published online 18 October 2007

SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB

Elmar Pruesse^{1,2}, Christian Quast^{1,3}, Katrin Knittel⁴, Bernhard M. Fuchs⁴, Wolfgang Ludwig⁵, Jörg Peplies⁶ and Frank Oliver Glöckner^{1,3,*}

¹Microbial Genomics Group, Max Planck Institute for Marine Microbiology, ²University Bremen, Center for Computing Technologies, D-28359, ³Jacobs University Bremen gGmbH, D-28759, ⁴Department of Molecular Ecology, Max Planck Institute for Marine Microbiology, D-28359 Bremen, ⁵Department for Microbiology, Technical University Munich, D-85354 Freising and ⁶Ribocon GmbH, D-28359 Bremen

- Useful databases:

- miRbase (<http://microrna.sanger.ac.uk/>)
- Rfam (<http://rfam.sanger.ac.uk/>)
- Silva (<http://www.arb-silva.de/>)
- GtRNAdb(<http://gtrnadb.ucsc.edu/>)



- Contains tRNA gene predictions made by the program tRNAscan-SE (Lowe & Eddy, Nucl Acids Res 25: 955-964, 1997) on complete or nearly complete genomes.
- All annotation is automated and has not been inspected for agreement with published literature.

Published online 4 November 2008

*Nucleic Acids Research, 2009, Vol. 37, Database issue D93-D97
doi:10.1093/nar/gkn787*

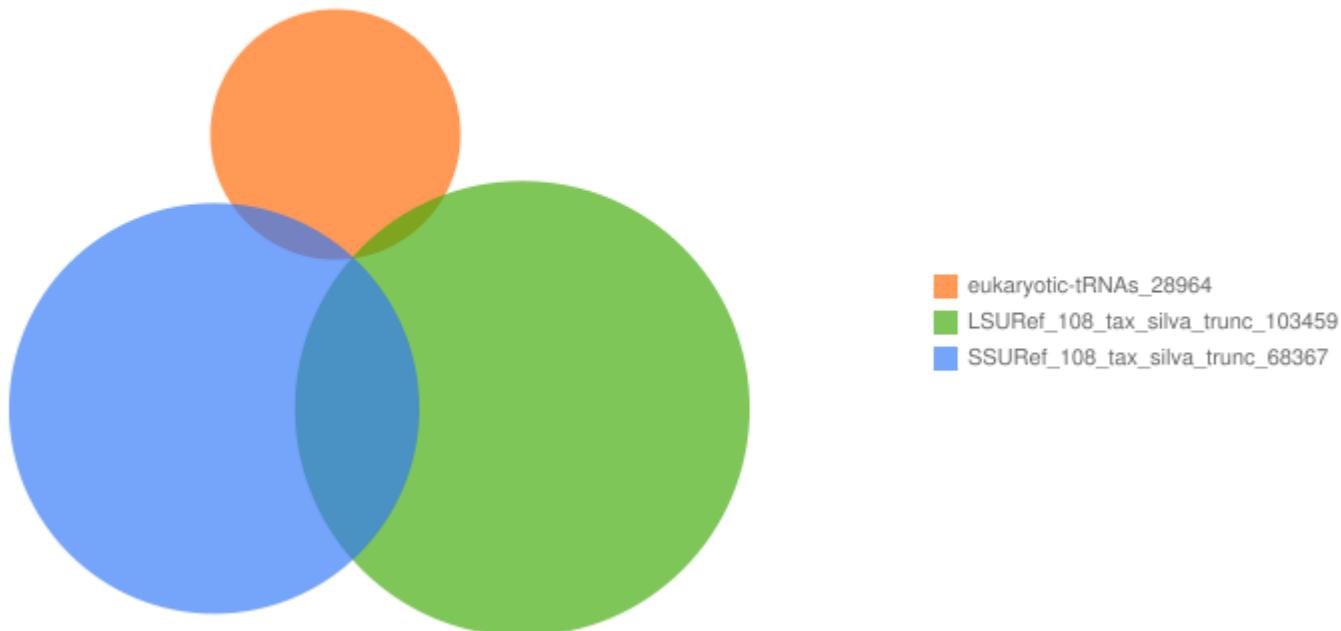
GtRNAdb: a database of transfer RNA genes detected in genomic sequence

Patricia P. Chan and Todd M. Lowe*

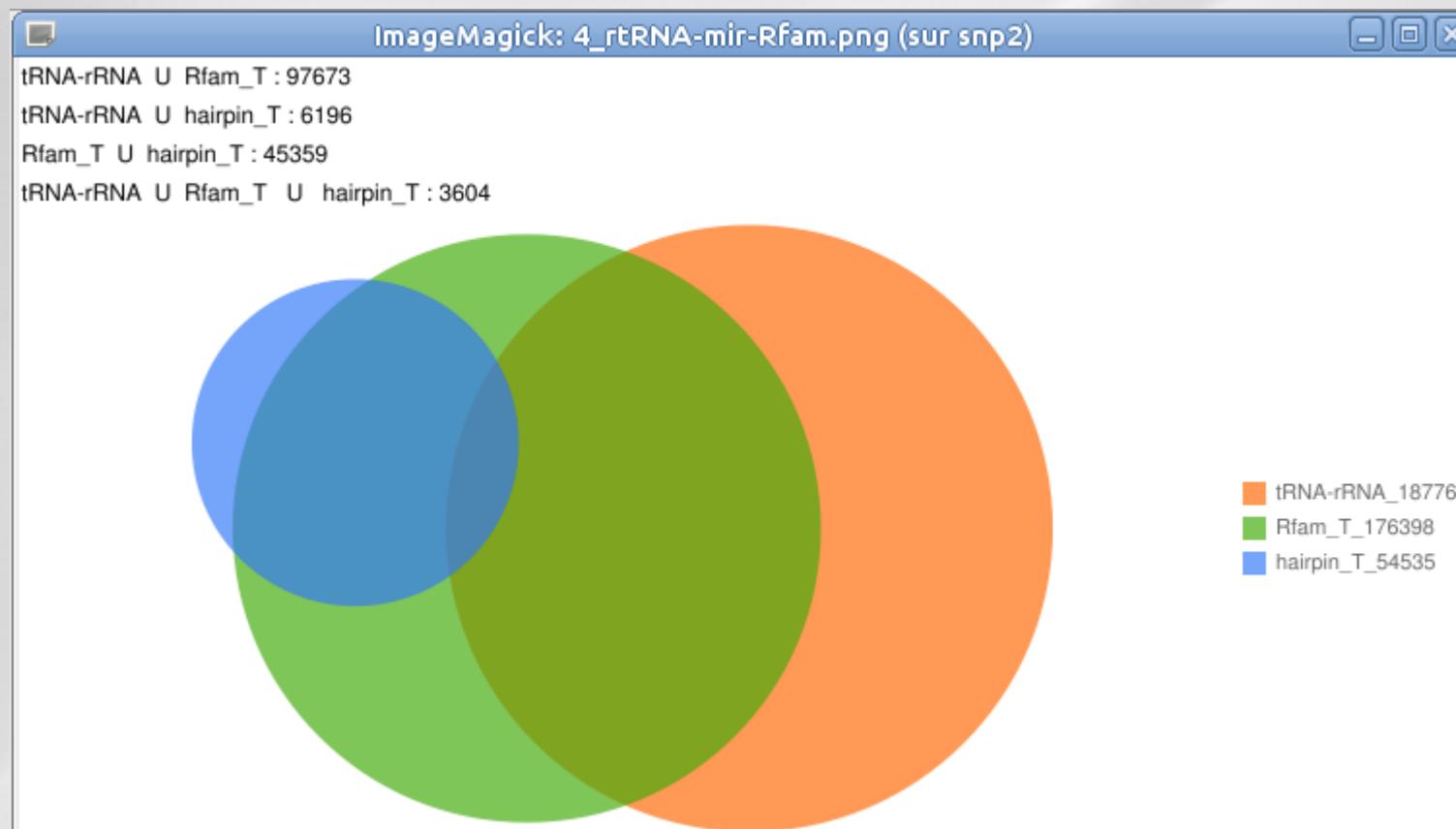
Department of Biomolecular Engineering, University of California, Santa Cruz, 1156 High Street, SOE-2, Santa Cruz, CA 95064, USA

- **Reads with multiple annotation**

```
eukaryotic-tRNAs U LSURef_108_tax_silva_trunc : 707
eukaryotic-tRNAs U SSURef_108_tax_silva_trunc : 1230
LSURef_108_tax_silva_trunc U SSURef_108_tax_silva_trunc : 11385
eukaryotic-tRNAs U LSURef_108_tax_silva_trunc U SSURef_108_tax_silva_trunc : 293
```



- **Reads with multiple annotation**



→ A lot of reads annotated with mirBase but also with tRNA and rRNA database

- rRNA present in miRBase

Mir-739 or 28S rRNA ?

Annotation

Annotation

occurrences

#seq	eukaryotic-tRNAs	hairpin_T	LSURef_108_tax Silva_trunc	Rfam_T	SSURef_108_tax Silva_trunc	SupportedBy	Total	s_1_ut21	s_1_ut2	s_1_ut4
							Search all columns:			
seq681297#1#189	0	oan-mir-20a-1	X54512.4749.8508	RF00051;mir-17;AAPN0128049.1/1987-2067	0	1	189	0	0	189
seq299078#2#304	0	mmu-mir-5105	V01270.3862.8647	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	2	304	165	0	0
seq610618#2#267	0	sha-mir-5105	V01270.3862.8647	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	2	267	102	0	0
seq1353575#4#218	0	mmu-mir-5105	U34342.1.3663	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	4	218	95	0	17
seq1353596#4#550	0	mmu-mir-5105	U34342.1.3663	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	4	550	161	0	183
seq2060361#3#113	0	mmu-mir-5105	U34342.1.3663	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	3	113	55	0	15
seq2060376#4#266	0	mmu-mir-5105	U34342.1.3663	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	4	266	97	3	56
seq1163251#5#342	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	5	342	96	2	116
seq1353595#5#239	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	5	239	57	4	111
seq1353600#5#759	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	5	759	170	29	247
seq2060374#4#113	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	4	113	25	0	62
seq401616#3#139	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	3	139	54	0	0
seq577112#4#524	0	mmu-mir-5105	U34341.1.3576	RF01960;SSU_rRNA_eukarya;AAYZ01438197.1/1-1685	0	4	524	146	0	203
seq1748431#4#548	0	cfa-mir-195	U34340.1.3432	RF00177;SSU_rRNA_bacteria;EU328070.1/1-1479	EU328070.1.1479	4	548	232	0	92
seq345104#4#102	0	gga-mir-1617	HQ856851.1.2611	RF00090;SNORA74;CAAE01008763.1/14090-14288	0	4	102	25	0	20
seq41650#5#523	0	sha-mir-716a	HQ856851.1.2611	RF00001;5S_rRNA;ABIM01036847.1/2163-2281	0	5	523	258	2	34
seq709529#5#160	0	hsa-mir-4792	GU372691.11134.15878	RF00100;7SK;AANN01516090.1/17881-17571	0	5	160	23	1	80
seq257457#2#119	0	sha-mir-716b	GQ424316.1.1993	RF00001;5S_rRNA;AARH01008767.1/1334-1421	0	2	119	0	0	106
seq718037#4#193	0	mmu-mir-5102	FP929060.89.2972	RF00028;Intron_gpf;EU352794.1/2419-2809	0	4	193	39	0	86
seq53378#5#144	0	mmu-mir-677	FP565809.564563.566970	RF01960;SSU_rRNA_eukarya;AAQR01407656.1/1-1561	AF198113.1.1740	5	144	43	3	56
seq1328312#4#393	0	ata-MIR172	FJ966040.1.2409	RF00100;7SK;AAQQ01276673.1/1502-1765	CABZ01109011.107.1605	4	393	155	24	0
seq1328326#4#142	0	ata-MIR172	FJ966040.1.2409	RF00306;snoZ178;AAZX01013617.1/1306-1470	CABZ01109011.107.1605	4	142	52	8	0
seq487403#4#645	0	ata-MIR172	FJ966040.1.2409	RF00306;snoZ178;AAZX01015218.1/4829-4668	U94741.1.2950	4	645	226	4	0
seq487443#4#169	0	sbi-MIR396c	FJ966040.1.2409	RF00100;7SK;AAKN02002849.1/102766-102498	CABZ01109011.107.1605	4	169	69	2	0
seq1328328#5#144	0	smo-MIR1082a	FJ966040.1.2409	RF00306;snoZ178;AC114644.10/51094-51230	CABZ01109011.107.1605	5	144	52	11	5
seq653494#4#168	0	mmu-mir-5102	FJ605292.1.3569	RF01960;SSU_rRNA_eukarya;CABB01000342.1/31007-29320	0	4	168	53	0	34
seq686909#5#164	0	r1cv-mir-rL1-8	FJ424422.1.2497	RF01960;SSU_rRNA_eukarya;BZ8748.1/1-1822	GQ352554.1.1846	5	164	6	4	140
seq1328311#5#316	0	ata-MIR172	FJ360703.1.2869	RF00009;RNaseP_nuc;AC102108.12/162476-162168	CABZ01109011.107.1605	5	316	80	24	6
seq667010#4#118	0	mmu-mir-5102	FJ040535.1.4142	RF00028;Intron_gpf;EU352794.1/2419-2809	0	4	118	42	0	8
seq1328321#4#323	0	osa-MIR408	EU921138.1.2387	RF00306;snoZ178;AAZX01015218.1/4829-4668	CABZ01109011.107.1605	4	323	91	23	0
seq487405#4#315	0	smo-MIR1082a	EU921138.1.2387	RF00306;snoZ178;AACO2015737.1/1625-1475	CABZ01109011.107.1605	4	315	124	3	0
seq1461535#5#1418	0	hsa-mir-4700	EU875589.109747.113671	RF00002;5_BS_rRNA;AJ270036.1/1-105	DM486508.4754.6504	5	1418	412	45	476
seq1861043#4#142	0	hsa-mir-4700	EU875589.109747.113671	RF00002;5_BS_rRNA;AF342795.1/144-297	AC211391.79568.81654	4	142	61	0	8

Exercice:

– Annotation