# small RNAseq data analysis
## miRNA detection
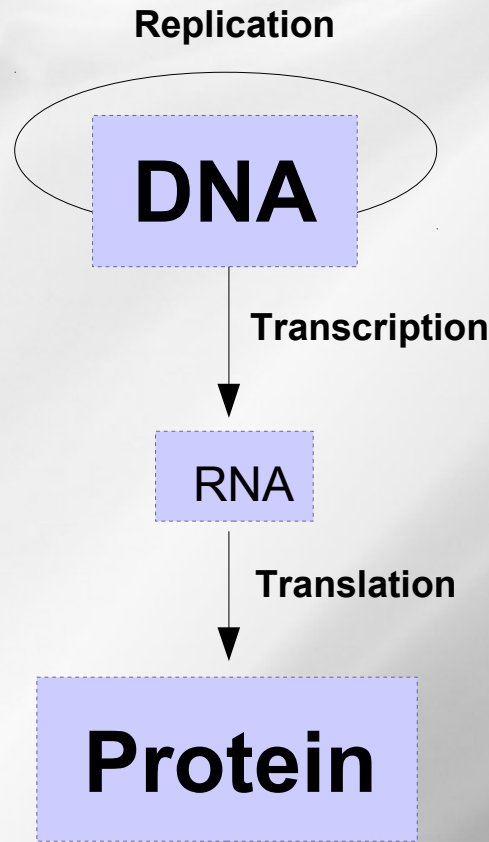
**P. Bardou, C. Gaspin, S. Maman, J. Mariette & O. Rué**

# Introduction ncRNA

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

  DNA structure discovery.

**Replication**

**DNA**

**Transcription**

RNA

**Translation**

**Protein**

**One gene = one function**

Adapted from Dujon B., Toulouse, 2009

- **Evolution of the dogma : 1970-1980**

  Genome analysis



Adapted from Dujon B., Toulouse, 2009

# Central dogma of molecular biology

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

  Genome analysis + Sequencing

**Replication**

DNA

**Reverse transcription**     **Transcription**

**Regulation Gene formation Epigenesis**

# RNA

**Splicing and edition events**

**Translation**

## ProteinS

# Many genes = one functionnel complex

Adapted from Dujon B., Toulouse, 2009

Plateforme Bioinformatique Midi-Pyrénées

- **An expending universe of RNA**

```
                          ┌─────────┐
                          │   RNA   │
                          └────┬────┘
            ┌──────────────────┴──────────────────┐
┌───────────────────────────┐           ┌──────────────────┐
│  mRNA <Riboswitches>       │           │  Non coding RNA  │
└───────────────────────────┘           └──────────────────┘
```

**« Regulatory » RNAs**

- miRNA (development...)
- siRNA (defense)
- piRNA (epigenetic and post-transcriptional gene silencing in spermatogenesis)

- sRNA (adaptative responses in bacteria)

**« Housekeeping » RNAs**

- telomerase RNA (replication)
- snRNA (maturation-splicing)
- snoRNA, gRNA (modification, editing)

**« Catalytic » RNAs**

- Hairpin ribozyme
- Hammerhead ribozyme
- …

→ **Multiple roles of RNA in genes regulation**

- **An expending universe of RNA**

```
                          ┌─────────┐
                          │   RNA   │
                          └────┬────┘
          ┌────────────────────┴───────────────────┐
┌──────────────────────────┐          ┌──────────────────────┐
│ mRNA <Riboswitches>      │          │   Non coding RNA     │
└──────────────────────────┘          └──────────────────────┘
```

**« Regulatory » RNAs**

- miRNA (development...)
- siRNA (defense)
- piRNA (epigenetic and post-transcriptional gene silencing in spermatogenesis)

- sRNA (adaptative responses in bacteria)

**« Housekeeping » RNAs**

- telomerase RNA (replication)
- snRNA (maturation-splicing)
- snoRNA, gRNA (modification, editing)

**« Catalytic » RNAs**

- Hairpin ribozyme
- Hammerhead ribozyme
- …

→ **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**
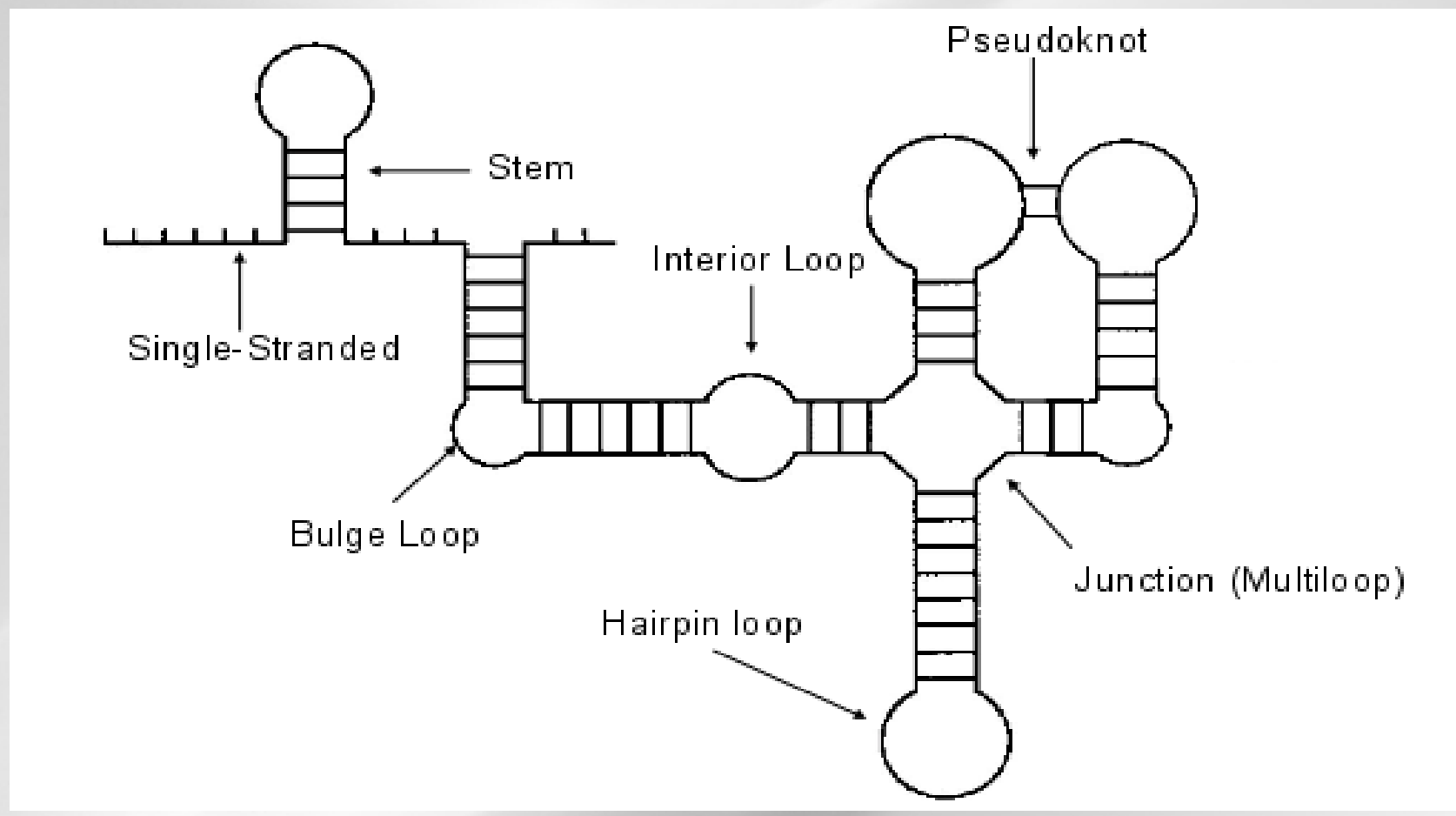
- **Sequence less conserved**

# 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

- **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 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 sRNAseq

- ## Discovery of lin-4 in C. elegans in 1993



```
5'GUUCCCUGAGACCUCAAGUG.UGAG   lin-4
3'CAAG.GACUC......UCGU-ACUC
  UAAG.GACUC.-.........ACUU
  CAAGGACUC.-...UUUAC-GCUC     lin-14
  UAAG.GACUC.-........U.ACUC   3'UTR
  CAAGGACUC.......CAU..CUU
  CAAG.GACU.-.....UGU.-UUC
  CA.GGACUC.-.........ACUC
```



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 Developm
Cambridge, Massachusetts 02138

**Summary**

*lin-4* is essential for the normal ter
diverse postembryonic developmen
elegans. *lin-4* acts by negatively regu
LIN-14 protein, creating a temporal d

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 pat-
terns of cell lineage appropriate for the L1 stage (Ambros
and Horvitz, 1984).
   *lin-14* controls these stage-specific cell lineages by gen-
erating a temporal gradient of Lin-14 nuclear protein (Lin



(He & Hannon, Nature reviews, 2004)

- ## 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. Small[1] and Eric N. Olson[1]

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

**THE JOURNAL OF IMMUNOLOGY**

**Small RNAs Guide Hematopoietic Differentiation and Function**

Francisco Navarro and Judy Lieberman

This information is current as of December 28, 2011

*J Immunol* 2010;184;5939-5947
doi:10.4049/jimmunol.0902567
http://www.jimmunol.org/content/184...

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

**Regulation of mouse stomach development and Barx1 expression by specific microRNAs**

Byeong-Moo Kim[1,2,*,†], Janghee Woo[1,3,†], Chryssa Kanellopoulou[4] and Ramesh A. Shivdasani[1,2,‡]

Developmental Cell 11, 441–450, October, 2006 ©2006 Elsevier Inc.  DOI 10.1016/j.devcel.2006.09.009

**The Diverse Functions of MicroRNAs in Animal Development and Disease**

Wigard P. Kloosterman[1] and Ronald H.A. Plasterk[1,2,*]

[1]Hubrecht Laboratory
Centre for Biomedical Genetics

Since then, several g... RNA-cloning strategies t... vertebrates and invertebra...

Leading Edge
**Review**

**ELSEVIER**

**miSSING LINKS: miRNAs and plant development**
Christine Hunter and R Scott Poethig

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 m... characterized mutations that are now ... components or targets of miRNA-med...

Addresses
Plant Science Institute, Department of Biok... Pennsylvania, Philadelphia, Pennsylvania 1...

PTGS and co-suppression, whereas siRNAs of 24–26 nt (long siRNAs) are associated with long-range transmission of silencing signals and methylation of corresponding genomic regions (Figure 1) [4]. The role of siRNAs in plant PTGS has been reviewed recently [5,6] and so is not discussed in detail here.

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

**Origin, Biogenesis, and Activity of Plant MicroRNAs**

Olivier Voinnet[1,*]
[1]Institut de Biologie Moléculaire des Plantes, CNRS UPR2357–Université de Strasbourg, 67084 Strasbou...
*Correspondence: olivier.voinnet@ibmp-ulp.u-strasbg.fr
DOI 10.1016/j.cell.2009.01.046

MicroRNAs (miRNAs) are key posttranscriptional regulators of eukaryotic g... use highly conserved as well as more recently evolved, species-specific m... array of biological processes. This Review discusses current advances in o... origin, biogenesis, and mode of action of plant miRNAs and draws compa... zoan counterparts.

Current Opinion in Genetics & Develo...
This review comes from a themed issue...
Pattern formation and developmental m...
Edited by Anne Ephrussi and Olivier Po...

0959-437X/$ – see front matter
© 2003 Elsevier Ltd. All rights reserved.

DOI 10.1016/S0959-437X(03)00081-9

**Review Article**

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

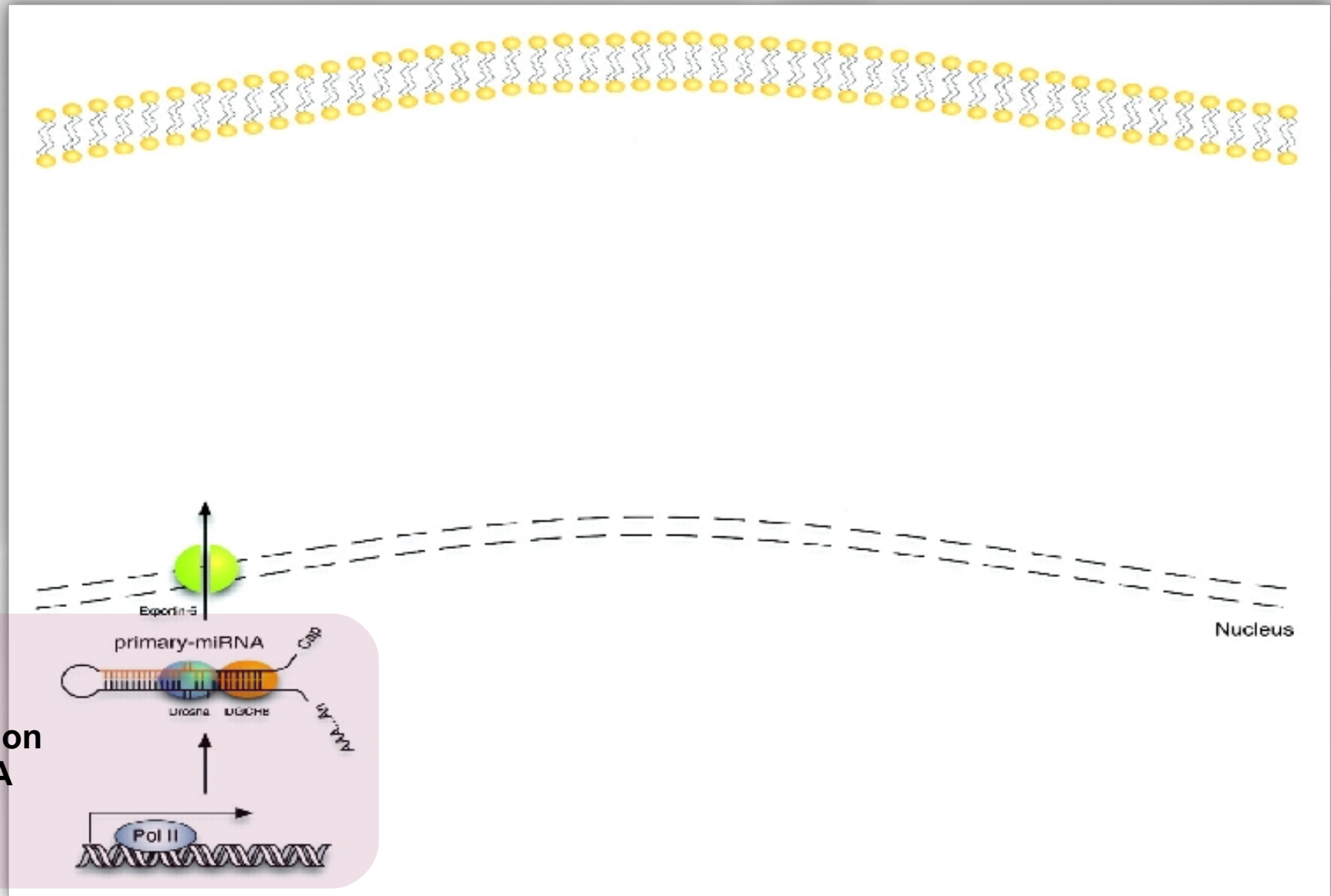Charlotte Delay[1,2] and Sébastien S. Hébert[1,2]

- # Animals

  - **Developmental timing (C. elegans):** lin-4, let-7

  - Neuronal left/right asymetry (C. elegans): Lys-6, mir-273

  - Programmed cell death/fat metabolism (D. melanogaster): mir-14

  - Notch signaling (D. malanogaster): mir-7

  - Brain morphogenesis (Zebrafish): mir-430

  - Myogeneses and cardiogenesis: mir-1, miR-181, miR-133

  - Insulin secretion: miR-375

  - …

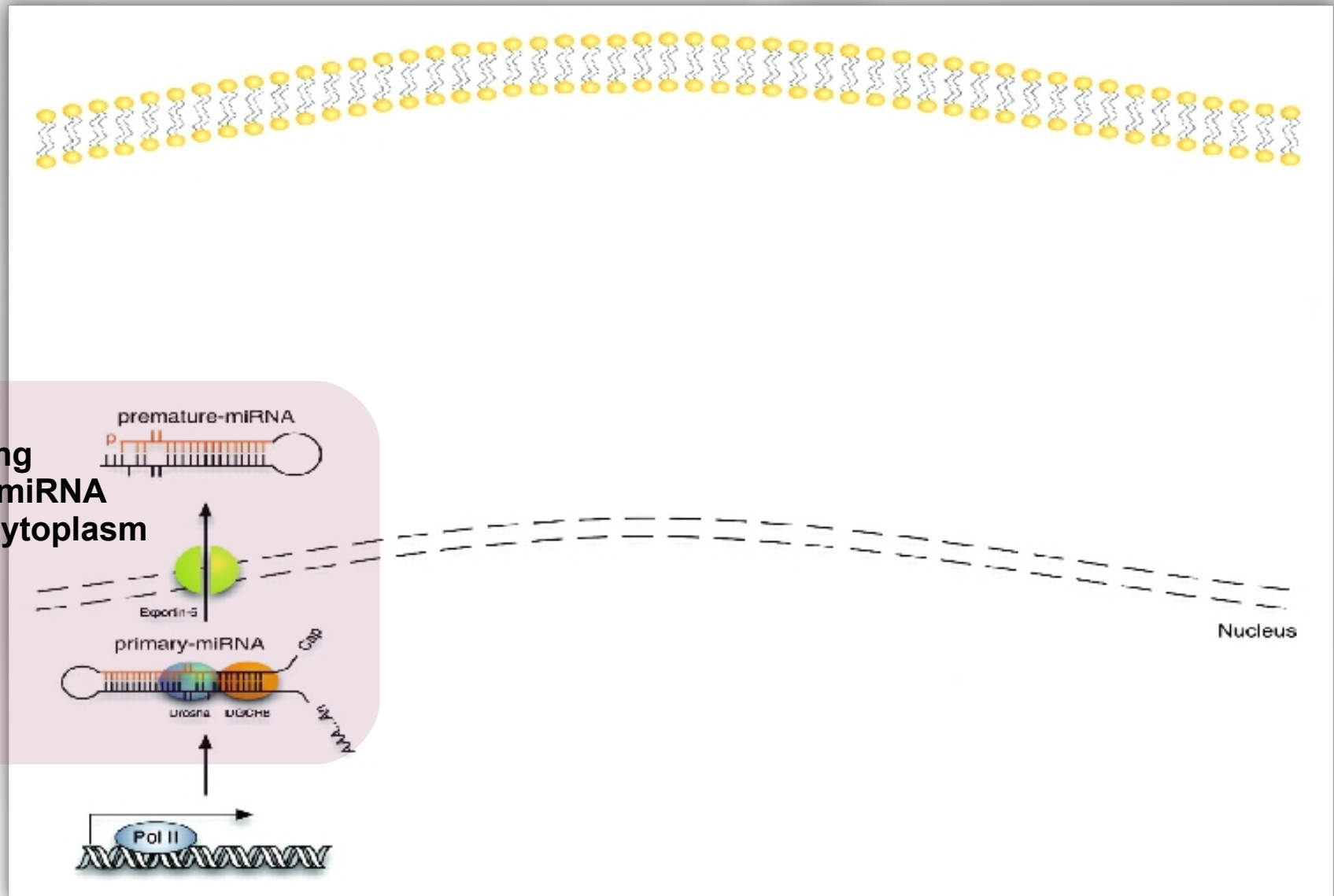  **1600 precursors in Human !!! (ref: miRBase, August 2012)**

- # Plants

  - **Floral timing and leaf development:** miR-156

  - Organ polarity, vascular and meristen development: mir-165, miR-166

  - Expression of auxin response genes: miR-160

  - ...

**Pol II transcription**
**Into a pri-miRNA**

**Drosha processing
one or more pre-miRNA
Exported in the cytoplasm**

One strand (miRNA)
Incorporated into
RISC

# The miRNA biogenesis

target mRNA translationally repressed

# The miRNA biogenesis



D. Bartel, Cell, 2004

a  Non-coding TU with intronic miRNA

miR-15a-16-1

DLEU2

b  Non-coding TU with exonic miRNA

miR-155

BIC

c  Coding TU with intronic miRNA

miR-25-93-106b

MCM7

d  Coding TU with exonic miRNA

miR-985

CACNG8

→ **Cluster organisation**

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

miRNA: plant vs animal

**Initial miRNA**

**Animal miRNA**

**Plant miRNA**

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

- RNAseq not suited for miRNA (protocol and size)



A. Poly-A selection, fragmentation and random priming

B. First and second strand synthesis

- 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

Plateforme Bioinformatique Midi-Pyrénées

- 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.

- **List of known miRNA**

- **List of new miRNA**

- **miRNA target(s)**

- **miRNA quantification**

- **Differential expression**

# small RNAseq 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

### Animal

Average : 87.8 nt

Average : 21.9 nt

### Plant

Average : 152.3 nt

Average : 21.3 nt

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

```
@D61655M1_171:2:1:1192:1017#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1192:1017#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:1202:1038#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1202:1038#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y[YXXcccccccc\cccc_aacccYUUVVOQ
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTTGCTTGAAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V]^BBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
+D61655M1_171:2:1:13819:1998#0/1
ggggggggggfgfggfg^ggggfggggeggggdgggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTTGTCAGACTTTTGTTTGGAGGTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
...
```

# Fastq format

```
@D61655M1_171:2:1:1192:1017#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1192:1017#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:1202:1038#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1202:1038#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y[YXXccccccc\cccc_aacccYUUVVOQ
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTTGCTTGAAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V]^BBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
+D61655M1_171:2:1:13819:1998#0/1
gggggggggfgfggfg^ggggfggggeggggdgggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTTGTCAGACTTTTGTTTGGAGGTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
. . .
```

## Line 1 starts with @

| Information | Meaning |
|---|---|
| D61655M1_171 | The unique instrument name |
| 2 | Flowcell lane45.156.426 |
| 1 | Tile number within the flox 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 (paired-end or mate-pair reads only) |

# Fastq format

```
@D61655M1_171:2:1:1192:1017#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1192:1017#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:1202:1038#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1202:1038#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y[YXXcccccccc\cccc_aacccYUUVVOQ
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTTGCTTGAAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V]^BBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
+D61655M1_171:2:1:13819:1998#0/1
gggggggggggfgfggfg^ggggfggggegggggdggggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTTGTCAGACTTTTGTTTGGAGGTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB

...
```

**Line 1** starts with @

| Information | Meaning |
|---|---|
| D61655M1_171 | The unique instrument name |
| 2 | Flowcell lane45.156.426 |
| 1 | Tile number within the flox 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 (paired-end or mate-pair reads only) |

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

# Fastq format

```
@D61655M1_171:2:1:1192:1017#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1192:1017#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:1202:1038#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1202:1038#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y[YXXccccccc\cccc_aacccYUUVVOQ
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTTGCTTGAAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V]^BBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
+D61655M1_171:2:1:13819:1998#0/1
gggggggggfgfggfg^ggggfggggegggggdgggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTTGTCAGACTTTTGTTTGGAGGTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB

...
```
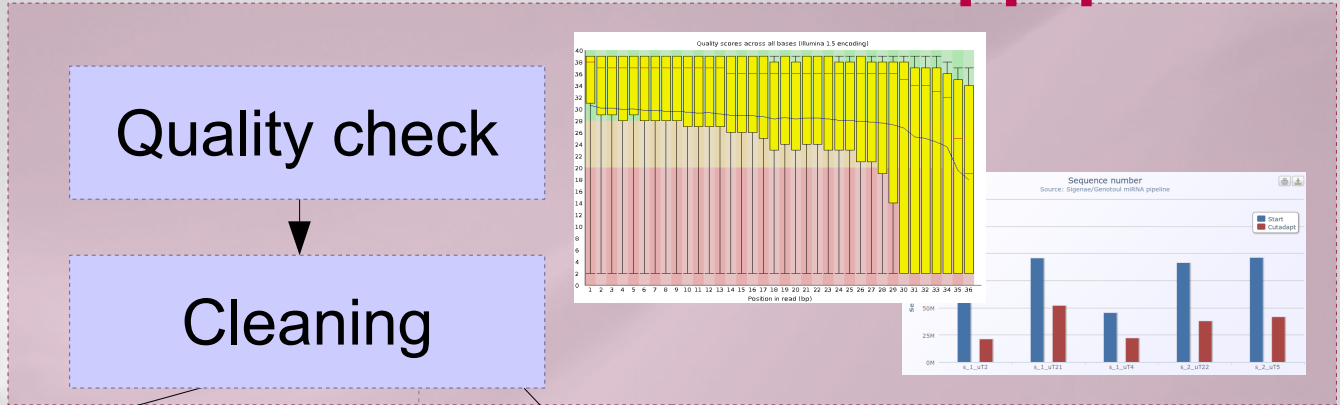
**Line 1** starts with @

| Information | Meaning |
|---|---|
| D61655M1_171 | The unique instrument name |
| 2 | Flowcell lane45.156.426 |
| 1 | Tile number within the flox 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 (paired-end or mate-pair reads only) |

**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

```
@D61655M1_171:2:1:1192:1017#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1192:1017#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:1202:1038#0/1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+D61655M1_171:2:1:1202:1038#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13360:1961#0/1
NTCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAA
+D61655M1_171:2:1:13360:1961#0/1
B[[[[Y[YXXcccccccc\cccc_aaccYUUVVOQ
@D61655M1_171:2:1:13406:1958#0/1
NGAGGTAGTAGATTGAATAGTTATCTCGTATGCCGT
+D61655M1_171:2:1:13406:1958#0/1
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13770:1993#0/1
GTCTCGTATGCCGGCTTTTGCTTGAAAAAAAAAGAA
+D61655M1_171:2:1:13770:1993#0/1
QV\^XQ\V]^BBBBBBBBBBBBBBBBBBBBBBBBBB
@D61655M1_171:2:1:13819:1998#0/1
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
+D61655M1_171:2:1:13819:1998#0/1
gggggggggfgfggfg^ggggfggggegggggdgggg
@D61655M1_171:2:1:2975:2145#0/1
TAGTTTGTCAGACTTTTGTTTGGAGGTCGTATGGCA
+D61655M1_171:2:1:2975:2145#0/1
^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
...
```

**Line 1** starts with @

| Information | Meaning |
|---|---|
| D61655M1_171 | The unique instrument name |
| 2 | Flowcell lane45.156.426 |
| 1 | Tile number within the flox 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 (paired-end or mate-pair reads only) |

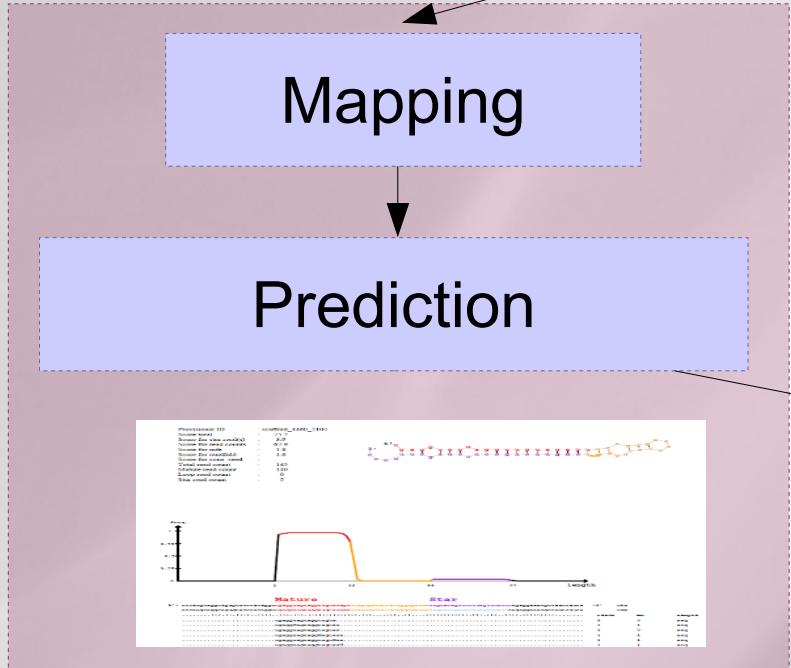**Line 2** Raw sequence of 36 nt (36 cycles in sequencing)

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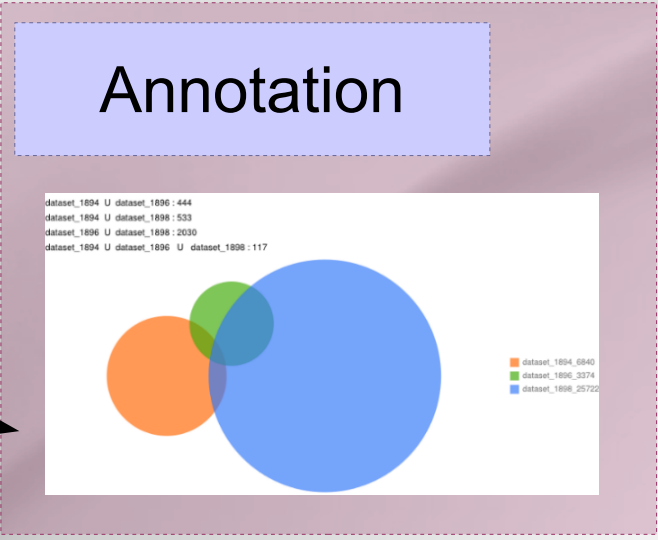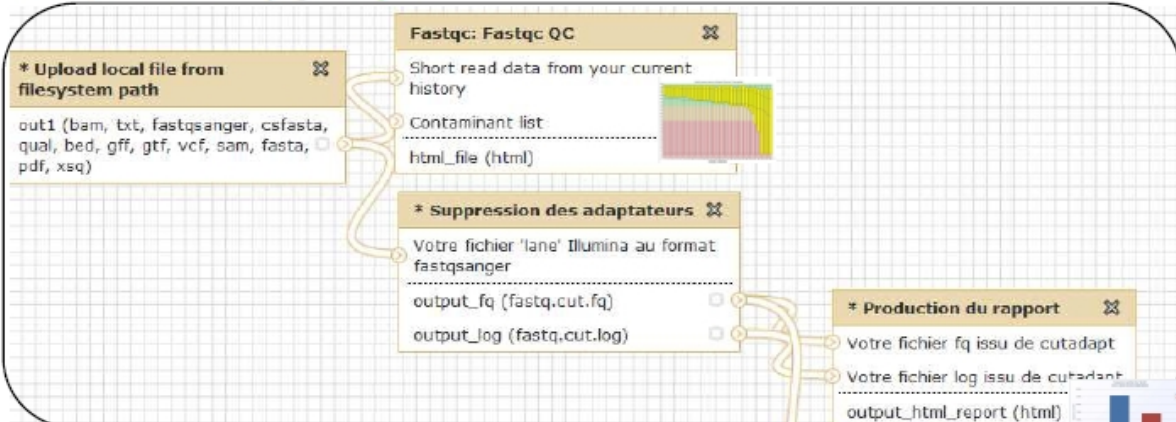
- # FastQC **(http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/)**

| Function | A quality control tool for high throughput sequence data. |
|---|---|
| Language | Java |
| Requirements | A suitable Java Runtime Environment<br><br>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
```

# Quality control

- **Per base quality**

- ## Sequences content in nucleotides

**Output reads**

```
>Adapteur
ATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
>UT1-10-28SrRNA
GCATGTTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40-piRNA ou tRNA
GCATTGGTGGTTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCTCGTATGCCGTCT
```

## Output reads

- Some sequences contain only adapters



```
>Adapteur
ATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
>UT1-10-28S rRNA
GCATGTTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40-piRNA ou tRNA
GCATTGGTGGTTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCTCGTATGCCGTCT
```

**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).

**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).

**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
ATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
>UT1-10-28S rRNA
GCATGTTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40-piRNA ou tRNA
GCATTGGTGGTTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCTCGTATGCCGTCT
```

## 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
ATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
>UT1-10-28S rRNA
GCATGTTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40-piRNA ou tRNA
GCATTGGTGGTTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCTCGTATGCCGTCT
```

**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
ATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAA
>UT1-10-28SrRNA
GCATGTTTGTGGAGAACCTGGTGCTAAATCACTCGT
>Poly-N
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
>UT1-40-piRNA ou tRNA
GCATTGGTGGTTCAGTGGTAGAATTCTCGCCATCTC
>UT1-2-mir21
TAGCTTATCAGACTGGTGTTGGCATCTCGTATGCCG
>UT1-3-mir143
TGAGATGAAGCACTGTAGCTATCTCGTATGCCGTCT
>UT1-30-mir143
TGAGATGAAGCACTGTAGCTCTCTCGTATGCCGTCT
```

## • **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<length<29).



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

- **56 % of reads discarded**



Sequence number
Source: Sigenae/Genotoul miRNA pipeline

- **Size in between 18bp:24bp**

  → **miRNA ?**



Length distribution
Source: Sigenae/Genotoul miRNA pipeline

# Exercices:

## – WF1:

- Quality control
- Cleaning

- ## 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
AAATGAATGATCTATGGACAGCAGC            51
AAATGAATGATCTATGGACAGCAGCA           82
AAATGAATGATCTATGGACAGCAGCAA          5
AAATGAATGATCTATGGACAGCAGCAAA         2
AAATGAATGATCTATGGACAGCAGCAAC         3
AAATGAATGATCTATGGACAGCAGCAAG         57
AAATGAATGATCTATGGACAGCAGCAG          2
AAATGAATGATCTATGGACAGCCGC            1
AAATGAATGATCTATGGACGGCAGCA           1
...
```

Length distribution after redundancy filter
Source: Sigenae/Genotoul miRNA pipeline

# Remove redundancy



miRNA ?    piRNA ?

Length distribution after redundancy filter
Source: Sigenae/Genotoul miRNA pipeline

- **More differencies between piRNAs than with miRNAs ?**

- 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**

- ...

- ## **Alignement of annotated reads**

- ## Alignement of annotated reads



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

- ## Alignement of all reads



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

# Exercices:

– **Mapping the reads with miRDeep2**

- **Using Bowtie for mapping**
- **miRDeep2-core for miRNA identification**

- 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

```
GAGAGTGGAGTGCAGCCAAGGATGACTTGCCGGAATTCACATATAGAGTGGAATGA
            CAGCCAAGGATGACTTGCCGG                    675
            CAGCCAAGGATGACTTGCCG                     26
              AGCCAAGGATGACTTGCCGG                    8
            CAGCCAAGGATGACTTGCCGGAA                   8
            CAGCCAAGGATGACTTG                         2
            CAGCCAAGGATGACTTGCCGGA                    2
            CAGCCAAGGATGACTTGC                        1
```

- Once the reads mapped

- Identify all contiguous read regions

- 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.

- Extend and fold read regions

- Extend and fold read regions

~ 100bp

- Extend and fold read regions

~ 100bp

- **Stable hairpin structure shifted by 2 bases**
- **miRNA > miRNA***

- Extend and fold read regions

- # Extend and fold read regions

**~ 100bp**



- **In the absence of reads corresponding to an expected miRNA*, additional checks on the structure are:**

  - 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 ?

genotoul bioinfo

Plateforme Bioinformatique Midi-Pyrénées

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 Hackenberg[1], Martin Sturm[2], David Langenberger[3,4],
Juan Manuel Falcón-Pérez[5] and Ana M. Aransay[1,*]

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

Published online 16 May 2010          *Nucleic Acids Research, 2010, Vol. 38, Web Ser*

## BMC Bioinformatics

BioMed Central

Open Access

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

Wei-Chi Wang[1], Feng-Mao Lin[1], Wen-Chi Chang[1,5], Kuan-Yu Lin[2,3],
Hsien-Da Huang*[1,4] and Na-Sheng Lin*[2,3]

Address: [1]Institute of Bioinforma
of Biotechnology, National Chen
Sinica, Nankang, Taipei 11529, T
Hsin-Chu 300, Taiwan, Republic
of China

Email: Wei-Chi Wang - canser.bi
Yu Lin - crisste@gate.sinica.edu.t

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

Genome **Biology**

## DSAP: deep-sequencing small RNA analys

Published online 12 September 2011          *Nucleic Acids Research, 2012, Vol. 40, No. 1  37–52*
doi:10.1093/nar/gkr688

**METHOD**                    Open Access

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

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

Marc R. Friedländer[1], Sebastian D. Mackow

*BIOINFORMATICS*   **APPLICATIONS NO**

*Sequence analysis*

### CPSS: a computational platform for the ana deep sequencing data

Yuanwei Zhang[1,†], Bo Xu[1,†], Yifan Yang[2], Rongjun Ban[3], H
Howard J. Cooke[1,4], Yu Xue[5,*] and Qinghua Shi[1,*]

[1]Hefei National Laboratory for Physical Sciences at Microscale and School of
and Technology of China, Hefei 230027, China, [2]Department of Statistics, Un
40506, USA, [3]Department of Computer Science & Technology, Nanjing Unive
Genetics Unit, IGMM, University of Edinburgh, Edinburgh EH4 2XU, UK, and [5]
Huazhong University of Science and Technology, Wuhan 430074, China

Associate Editor: Ivo Hofacker

NATURE BIOTECHNOLOGY VOLUME 26   NUMBER 4   APRIL 2008

## Discovering microRNAs from deep sequencing data using miRDeep

Marc R Friedländer[1], Wei Chen[2], Catherine Adamidi[1], Jonas Maaskola[1], Ralf Einspanier[3], Signe Knespel[1] &
Nikolaus Rajewsky[1]

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
DOI 10.1007/s11103-012-9885-2

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

*NOTE*

Advance Access publication August 27, 2010

ep sequencing analysis

ov[2], Gideon Dror[2], Eran Halperin[3,4]

dicine, Tel Aviv University, [2]The Academic
ce Institute, Berkeley, CA, USA and [4]School
iotechnology, George Wise Faculty of Life

## *shortran*: A pipeline for small RNA-seq data analysis

Vikas Gupta[1,2], Katharina Markmann[1], Christian N. S. Pedersen[2], Jens Stougaard[1] a
Andersen[1*]

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

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

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

- **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 |

- Reads pre-processing

  – Adaptators trimming

  – Redundancy

  – Repeats

  – Other ncRNA

  – Size of the mature miRNA (min/max)

- Mapping

  – Size and region of the read

  – Number of locations

    • Considered

    • Reported

  – Error(s) consideration in mapping

  – Quality of the read

- Precursor identification
  - Length and bounds of the theoretical sequence (folding)
  - Alignment of the read against known miRNA

- Post processing step: assessment of the potential miRNA
  - Different 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...)

NATURE BIOTECHNOLOGY VOLUME 26 NUMBER 4 APRIL 2008

# Discovering microRNAs from deep sequencing data using miRDeep

Marc R Friedländer[1], Wei Chen[2], Catherine Adamidi[1], Jonas Maaskola[1], Ralf Einspanier[3], Signe Knespel[1] & Nikolaus Rajewsky[1]

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. Friedländer[1], Sebastian D. Mackowiak[1], Na Li[2], Wei Chen[2] and Nikolaus Rajewsky[1,*]

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

Plateforme Bioinformatique Midi-Pyrénées

## Three modules

- **MiRDeep2**
- **Mapper**
- **Quantifier**

## 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

0 mismatch
In the seed

2 mismatches
after the seed

The size of the
Seed is 18nt

Do not map
more than 5

Order alignments
from best to worse

## The miRDeep2 module

- Scan of both strands from 5' to 3'

## The miRDeep2 module

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

## The miRDeep2 module

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

## The **miRDeep2** module

- Scan of both strands from 5' to 3
  - Search the best stack of reads (heigh 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 (RNAfold + Randfold)
    - Unbifurcated hairpins
  - Score the candidates
    - Valid alignment of reads on the precursor
    - 60% of nt in the mature part paired

## The **Quantifier** module

Identifies and quantifies known mature miRNA given

– Know mature miRNA

– Know miRNA precursors

Use Bowtie for miRNA/reads alignment

# Exercice:

- Back to miRdeep2-core results

- ## 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.



D152–D157   Nucleic Acids Research, 2011, Vol. 39, Database issue       Published online 30 October 2010
doi:10.1093/nar/gkq1027

# 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

## Rfam: updates to the RNA families database

Paul P. Gardner[1,*], Jennifer Daub[1], John G. Tate[1], Eric P. Nawrocki[2],
Diana L. Kolbe[2], Stinus Lindgreen[3], Adam C. Wilkinson[1], Robert D. Finn[1],
Sam Griffiths-Jones[4], Sean R. Eddy[2] and Alex Bateman[1]

[1]Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SA, UK, [2]Howard Hughes Medical Institute, Janelia Farm Research Campus, Ashburn, Virginia, USA, [3]Center for Bioinformatics, Department of Biology, University of Copenhagen, Ole Maaloes Vej 5, DK-2200 Copenhagen N, Denmark and [4]Faculty 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*  Published online 18 October 2007
*doi:10.1093/nar/gkm864*

## 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 Knittel[4], Bernhard M. Fuchs[4], Wolfgang Ludwig[5], Jörg Peplies[6] and Frank Oliver Glöckner[1,3,*]

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

**GtRNAdb: a database of transfer RNA genes detected in genomic sequence**
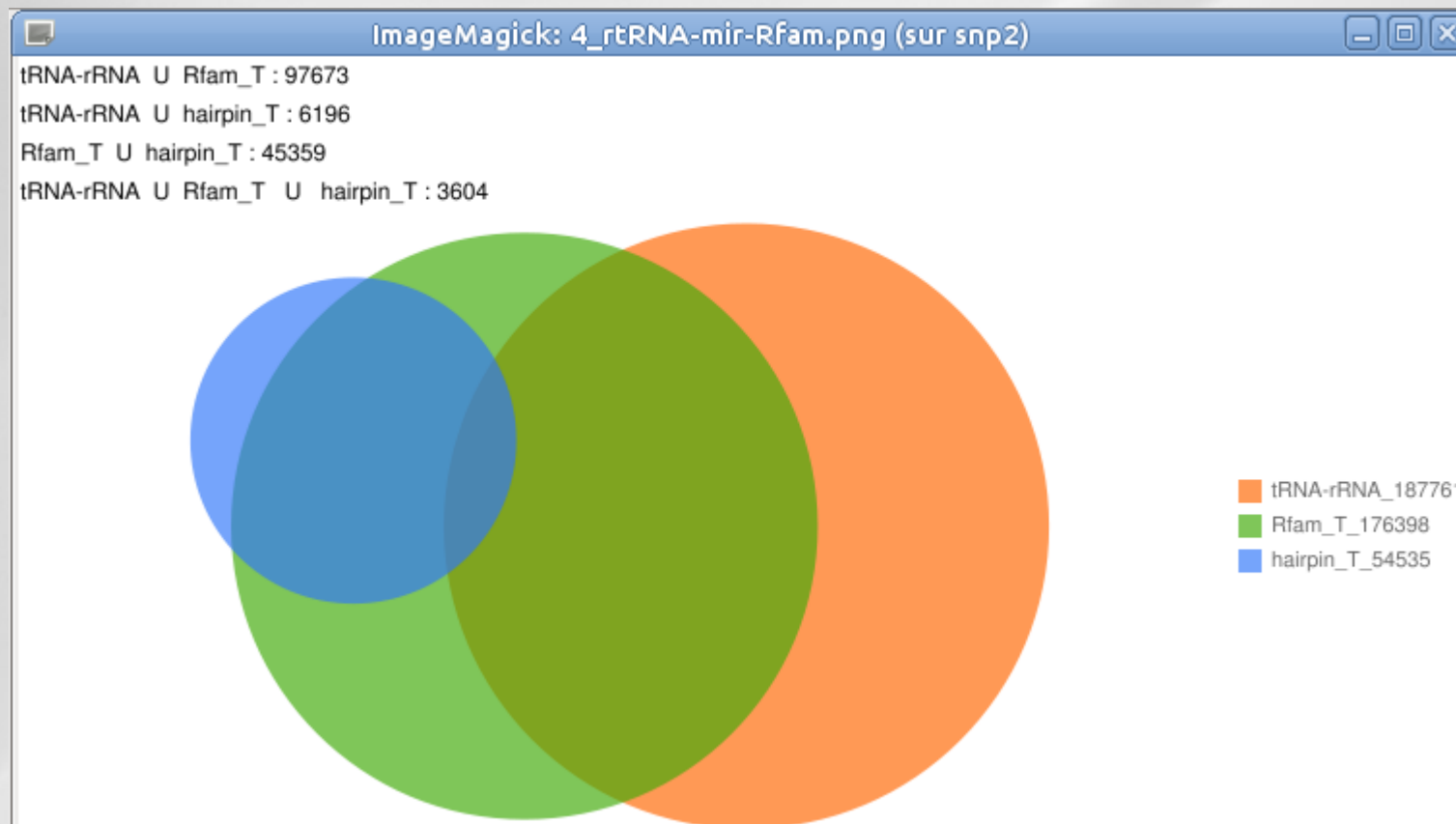
Patricia P. Chan and Todd M. Lowe*

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

- **Reads with multiple annotation**

- **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          occurences

| #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 |
|---|---|---|---|---|---|---|---|---|---|---|
| seq681297#1#189 | 0 | oan-mir-20a-1 | X54512.4749.8508 | RF00051;mir-17;AAPN01282042.1/1987-2067 | 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_gpI;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 | rlcv-mir-rL1-8 | FJ424422.1.2497 | RF01960;SSU_rRNA_eukarya;Z83748.1/1-1822 | 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_gpI;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;AASC02015737.1/1625-1475 | CABZ01109011.107.1605 | 4 | 315 | 124 | 3 | 0 |
| seq1461535#5#1418 | 0 | hsa-mir-4700 | EU875589.109747.113671 | RF00002;5_8S_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_8S_rRNA;AF342795.1/144-297 | AC211391.79568.81654 | 4 | 142 | 61 | 0 | 8 |

Show 100 entries        Search all columns:

# Exercice:

 − **Annotation**