# Module 2 : RNAseq bioinformatic
## from sequence to expression level
# smallRNA & isoform (eukaryote)

**Claire Toffano-Nioche**

I2BC - Institut de Biologie Intégrative de la Cellule, Orsay

# RNAseq analyses

RNAseq: From sequence data (reads) to expression level (count)

Classical analyses of RNA-Seq:
- – Check quality, Trimming
- – Mapping / counts
- – Assembly

Others usages of RNA-Seq:
- – smallRNA study
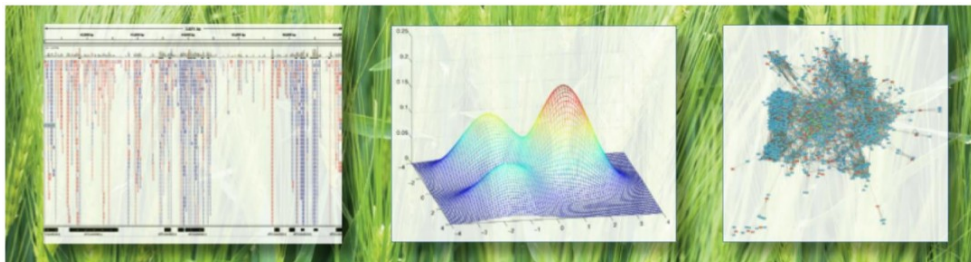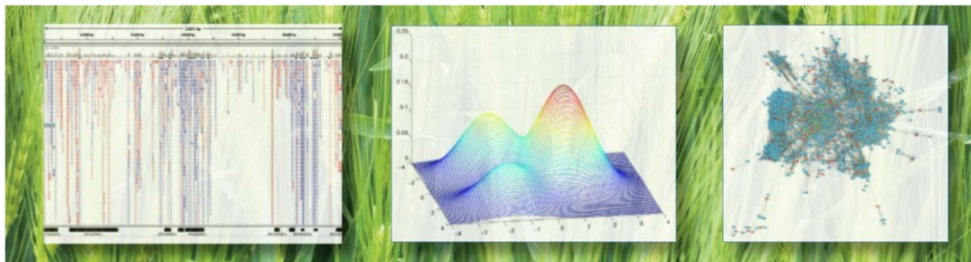- – expression at isoform level

# RNAseq analyses

RNAseq: From sequence data (reads) to expression level (count)
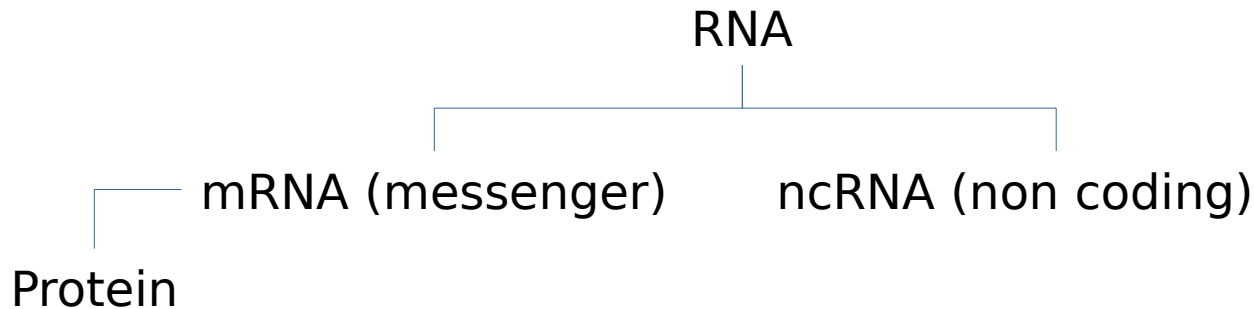
Classical analyses of RNA-Seq:
- Check quality, Trimming
- Mapping / counts
- Assembly

Others usages of RNA-Seq:
- **smallRNA study**
- expression at isoform level



SPS
SACLAY PLANT SCIENCES

# RNA world

RNA

mRNA (messenger)     ncRNA (non coding)

Protein

Not predicted by gene prediction tools
- No specific signal (start, stop, splicing sites...)
- Multiple location (intergenic, intronic, coding, antisense)
- Variable size
- No strong sequence conservation in general

Function related to structure : ncRNA of the same family have a conserved structure
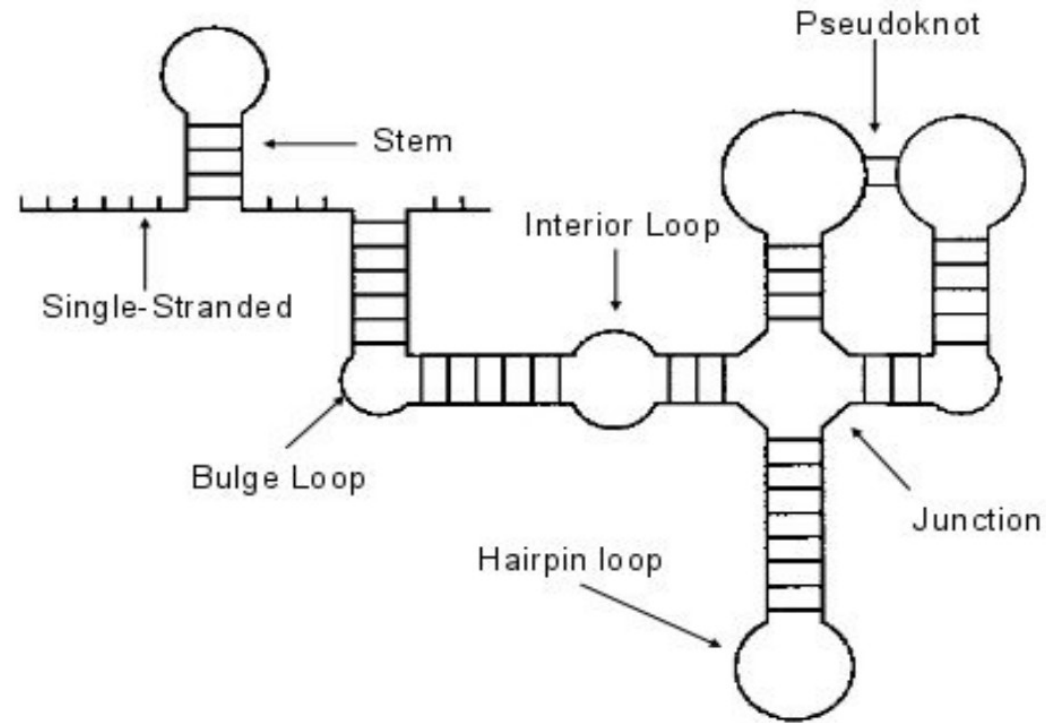
# RNA folding

Folding = Secondary structure

RNA folds on itself by base pairing :
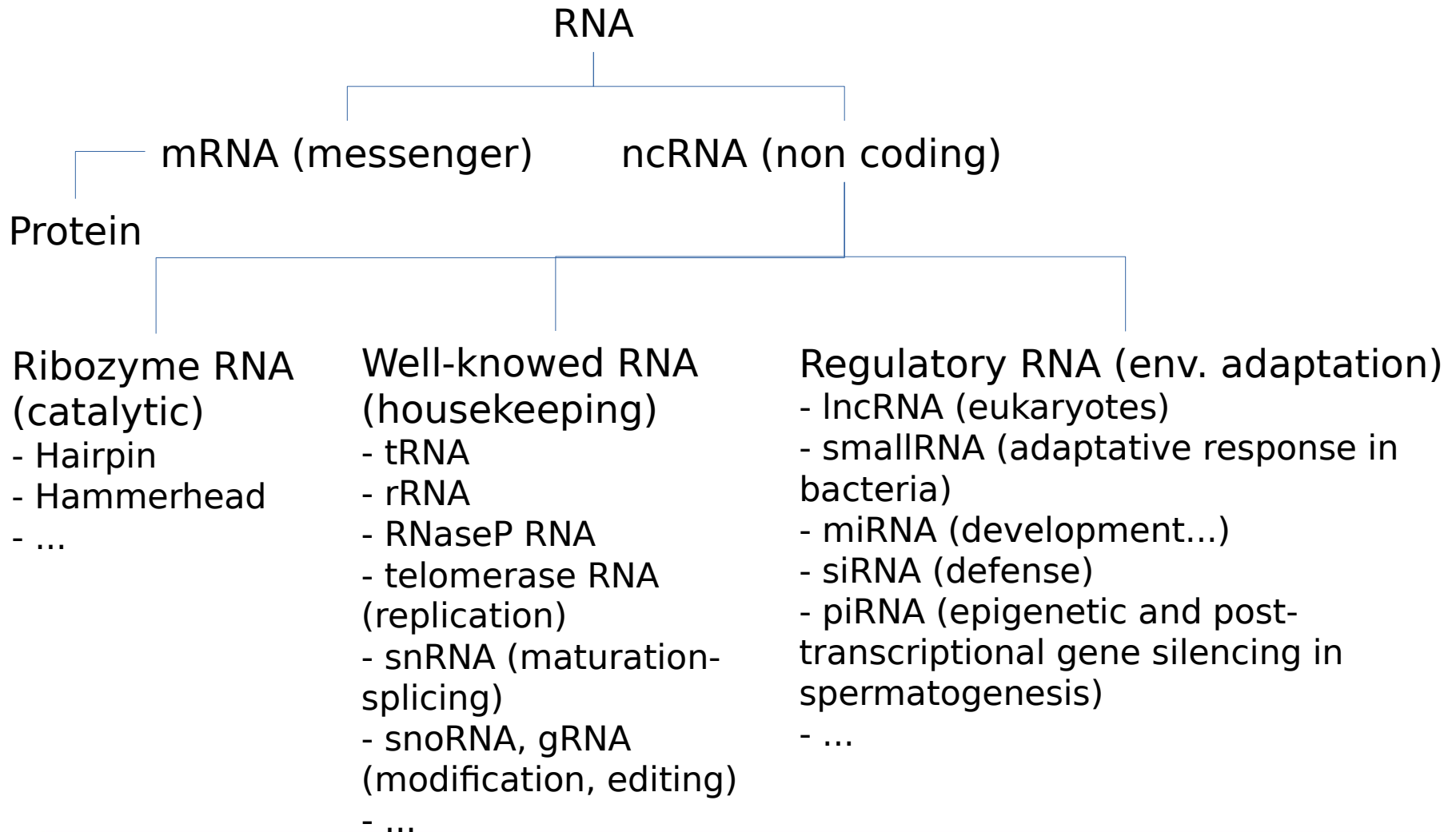
- A with U
- C with G
- sometimes G with U

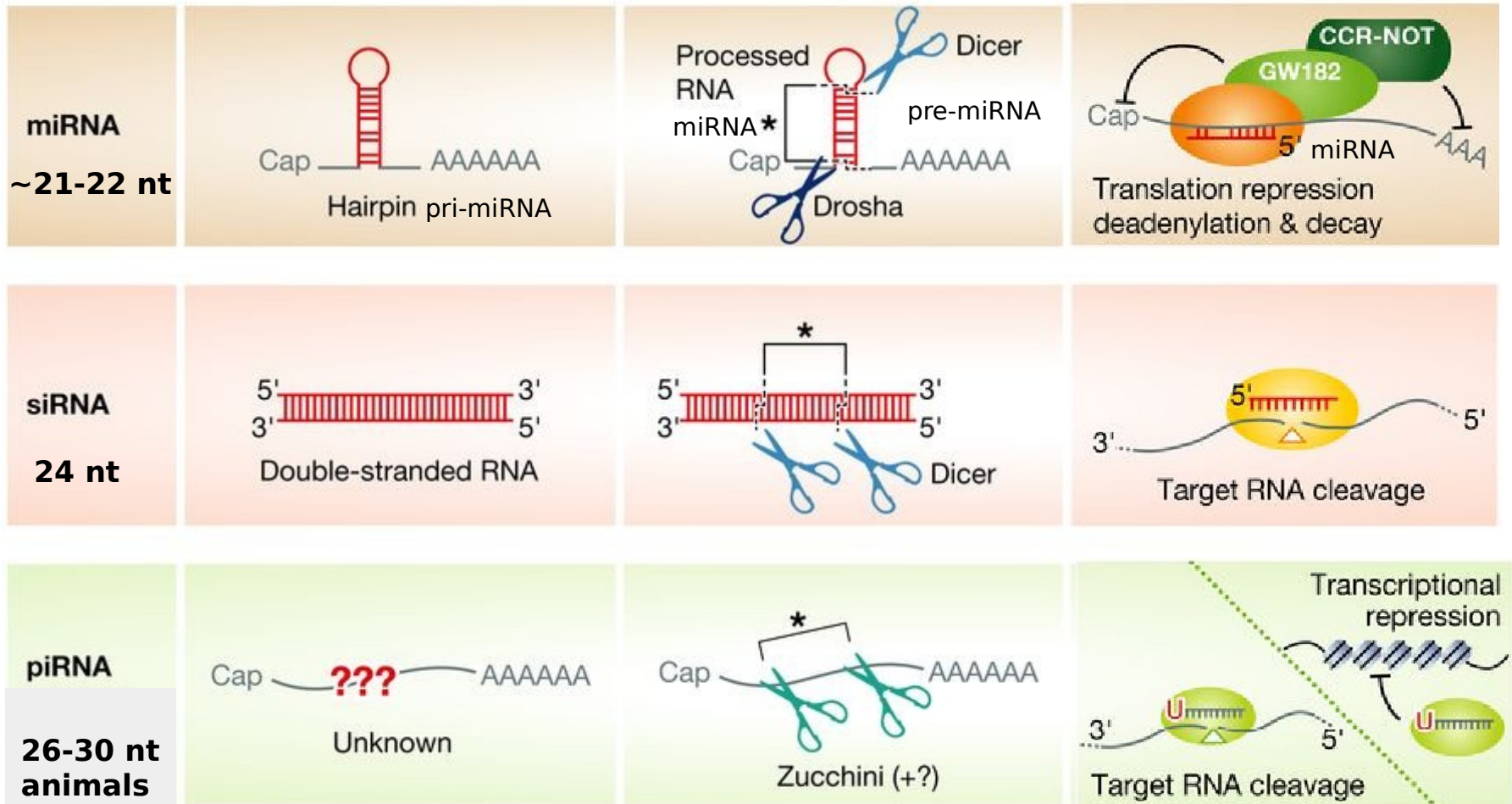=> more combinatorial possibilities of folding than DNA

RNA folds:
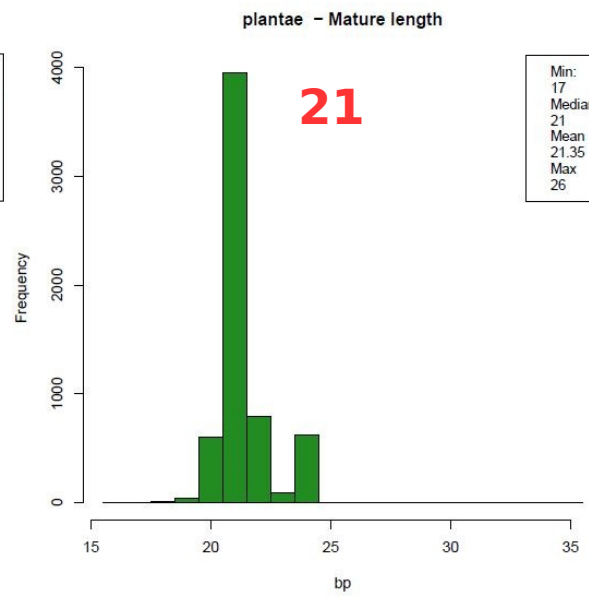- – on itself
- – with another: RNA duplex

# RNA world

RNA
- mRNA (messenger)
- ncRNA (non coding)

Protein

**Ribozyme RNA (catalytic)**
- Hairpin
- Hammerhead
- ...

**Well-knowed RNA (housekeeping)**
- tRNA
- rRNA
- RNaseP RNA
- telomerase RNA (replication)
- snRNA (maturation-splicing)
- snoRNA, gRNA (modification, editing)
- ...

**Regulatory RNA (env. adaptation)**
- lncRNA (eukaryotes)
- smallRNA (adaptative response in bacteria)
- miRNA (development...)
- siRNA (defense)
- piRNA (epigenetic and post-transcriptional gene silencing in spermatogenesis)
- ...

# Eukaryotic regRNA

Hirose T, Mishima Y, Tomari Y. (2014): Elements and machinery of non-coding RNAs: toward their taxonomy. EMBO Rep. 2014 May;15(5):489-507
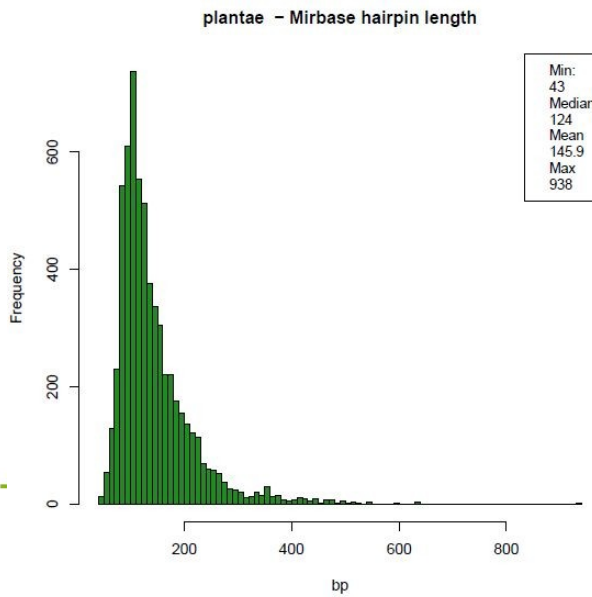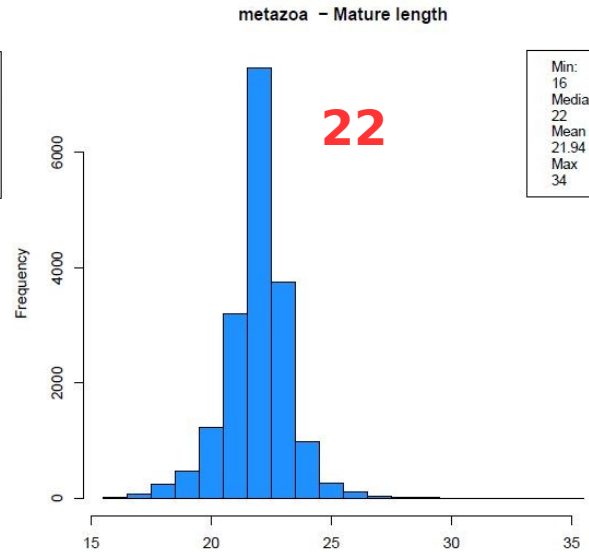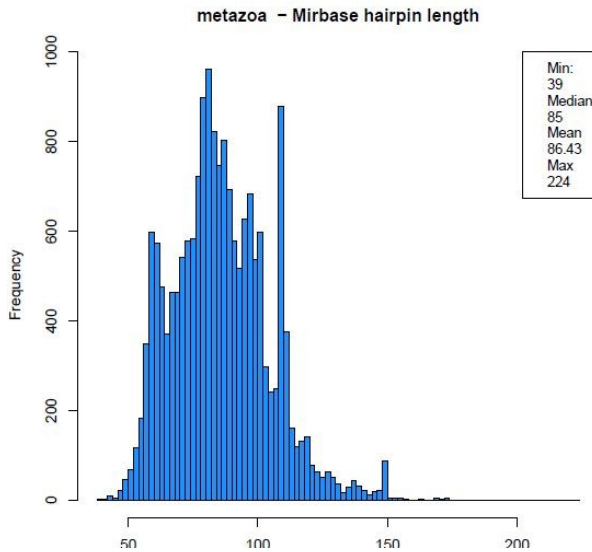
# Size maters



Griffiths-Jones S, Saini HK, van Dongen S, Enright AJ. miRBase: tools for microRNA genomics. Nucleic Acids Res. 2008 Jan;36(Database issue):D154-8

8

# smallRNA & NGS

Classical RNAseq not suited for smallRNA (protocol and size)

> smallRNAs lack a common sequence (e.g. a poly(A) tail) that can be used for selective enrichment or as a universal primer-binding site for reverse transcription

Strategies to enrich RNA sample in smallRNA (*):

- smallRNA cloning, deep sequencing : significant biases are introduced during small RNA cDNA library preparation (often more than 3 orders of magnitude for individual miRNAs ; one major source: RNA ligation)

- RNA Immuno-Precipitation (RIP-seq)

- Total RNA extraction + size selection

# smallRNAseq pipeline

Experimental design

Sequencing

Module 1 : High Troughput Sequencing
Technologies, available and futures...

Quality check

Cleaning

Module 2 : Bioinformatic of RNAseq

size selection

*w.o. reference*

*with reference*

Mapping

Prediction → Annotation

Module 3 : Normalization
& differential analysis of
RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

FastQC

Quality check

see previous talk

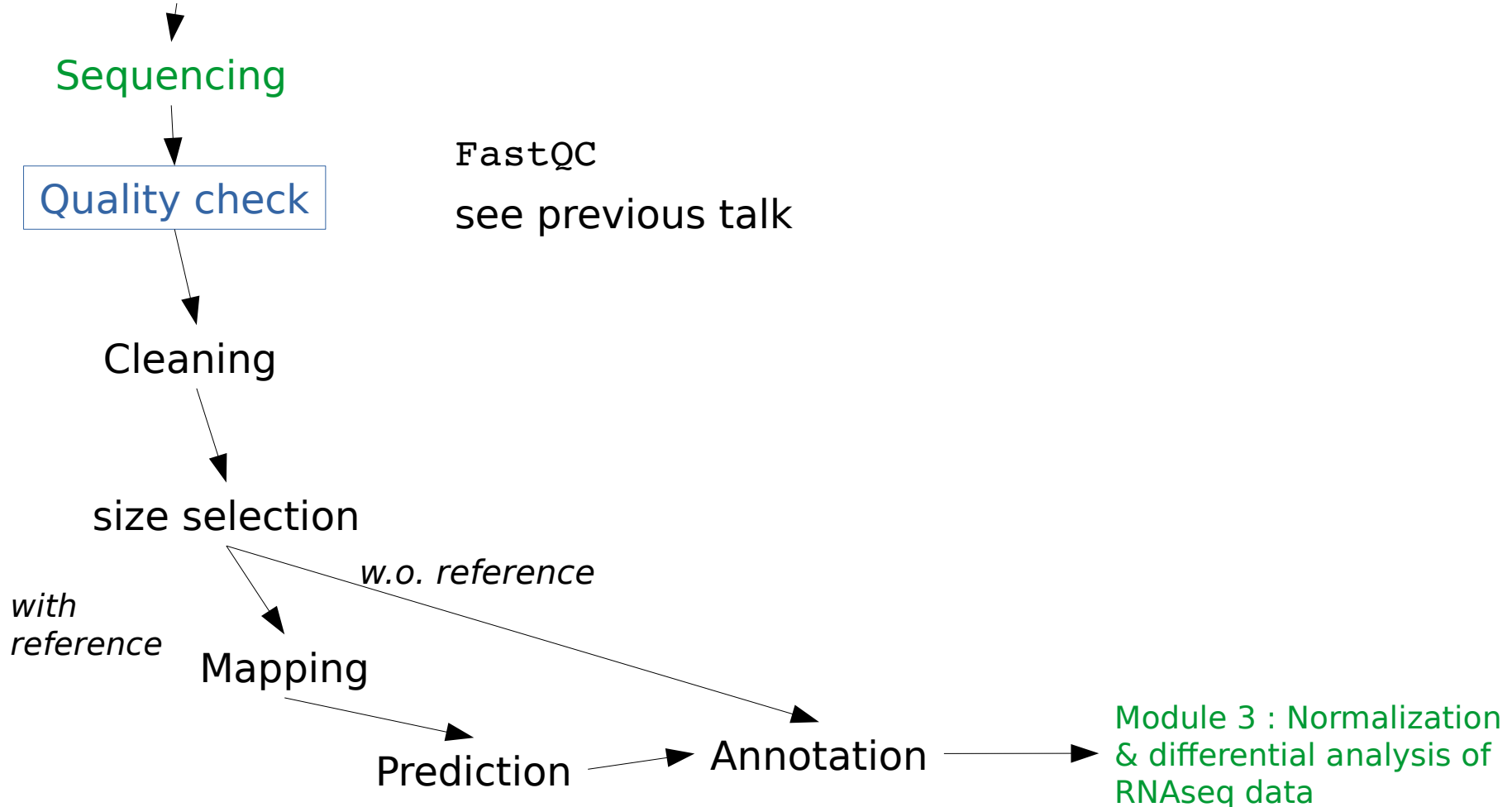Cleaning

size selection

*w.o. reference*

*with reference*

Mapping

Prediction

Annotation

Module 3 : Normalization & differential analysis of RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

*with reference*

Mapping

*w.o. reference*

Prediction

Annotation

Read size > regRNA size
PCRprimer sequences & adapter are included in read

## ❌ Overrepresented sequences

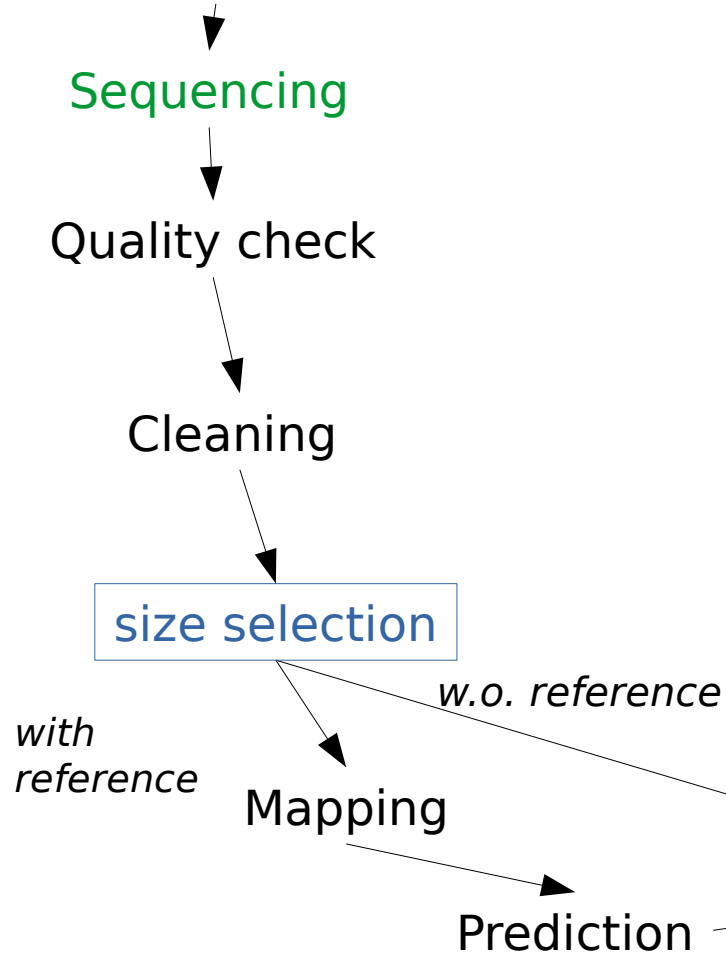| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GGGGATGTAGCTCAGAAGATCGGAAGAGCACACGTCTGAACTCCAGTCAC | 3865 | 1.546 | Illumina Multiplexing PCR Primer 2.01 (100% over 34bp) |
| GCGTCTGTAGTCCAACGGTTAGGATAATTGCAGATCGGAAGAGCACACGT | 3021 | 1.2084 | No Hit |
| GGGGATGTAGCTCAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCA | 2205 | 0.882 | TruSeq Adapter, Index 7 (100% over 35bp) |
| AGATCGGAAGAGCACACGTCTGAACTCCAGTCACCAGATCATCTCGTATG | 2047 | 0.8188000000000001 | TruSeq Adapter, Index 7 (100% over 49bp) |
| AGGGCTATAGCTAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCAGA | 1478 | 0.5912 | TruSeq Adapter, Index 7 (100% over 37bp) |
| CTAACAGACCGGTAGACTTGAACAGATCGGAAGAGCACACGTCTGAACTC | 1222 | 0.4888 | Illumina Multiplexing PCR Primer 2.01 (100% over 27bp) |
| CTTGAACAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCAGATCATC | 1155 | 0.462 | TruSeq Adapter, Index 7 (100% over 42bp) |

=> trimming step : remove adapter in 3', no quality trimming, remove rRNA, tRNA
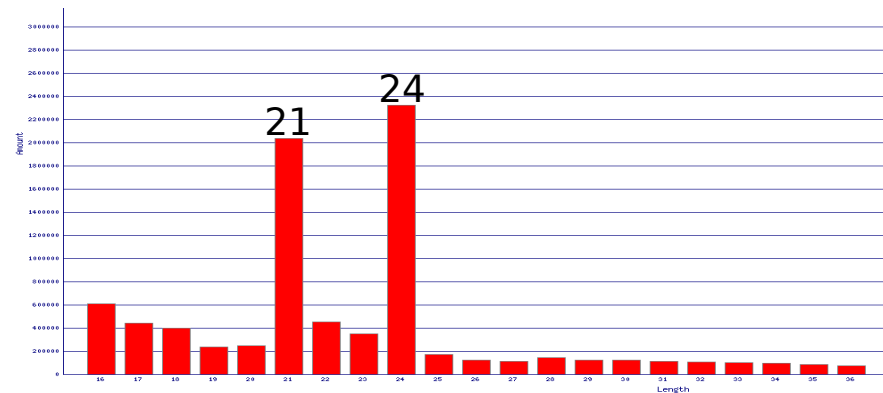
`trimgalore, cutadapt, trimmomatic, …`

Module 3 : Normalization & differential analysis of RNAseq data

12

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

*with reference*

*w.o. reference*

Mapping

Prediction

Annotation

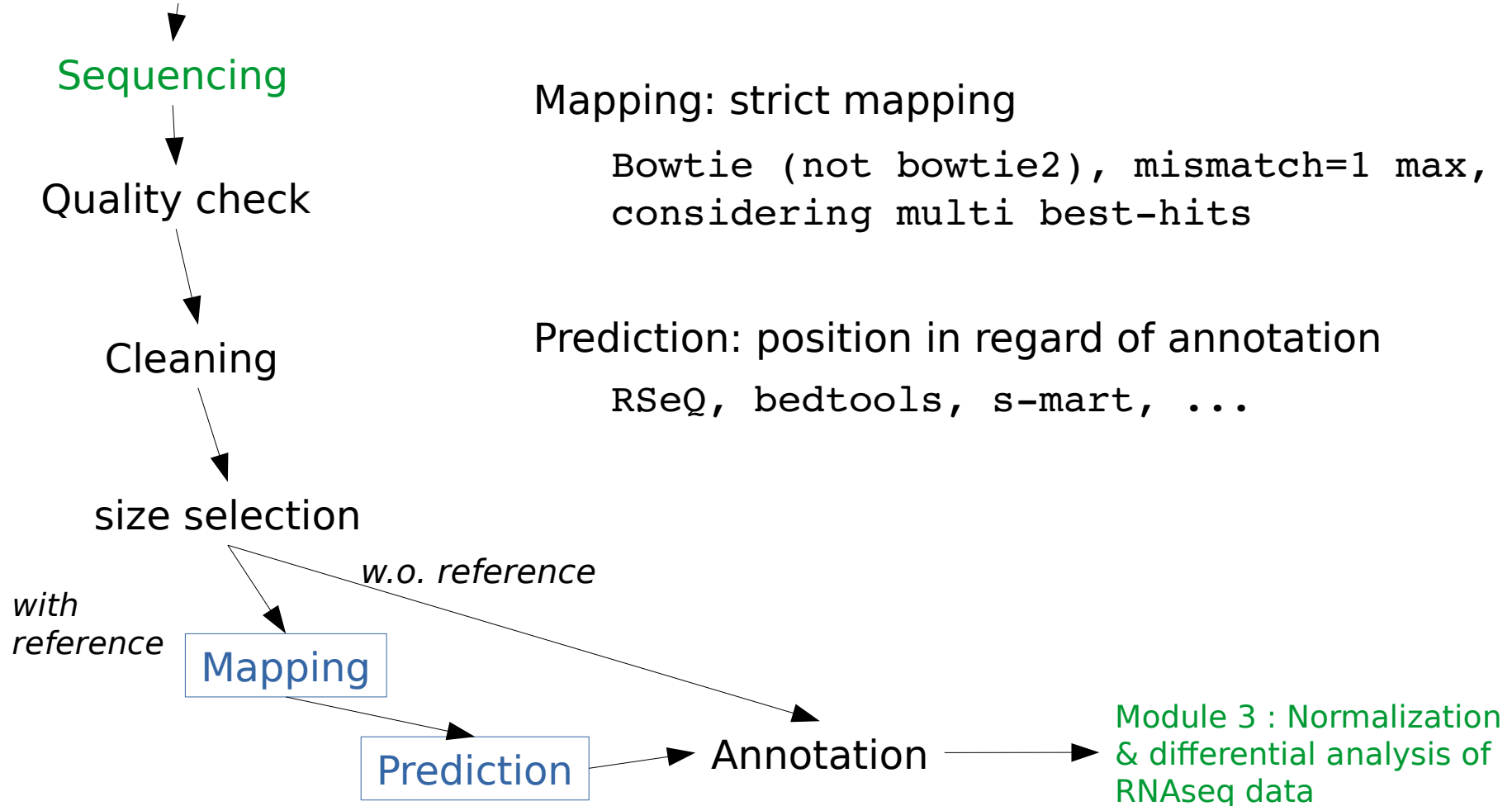Check the distribution of length reads
Plant: miRNA=21nt, siRNA=24nt



detect redonduncy

`fasta_clipping_histogram,`
`duplicate fastx_collapser`
from fastx-toolkit

Module 3 : Normalization
& differential analysis of
RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

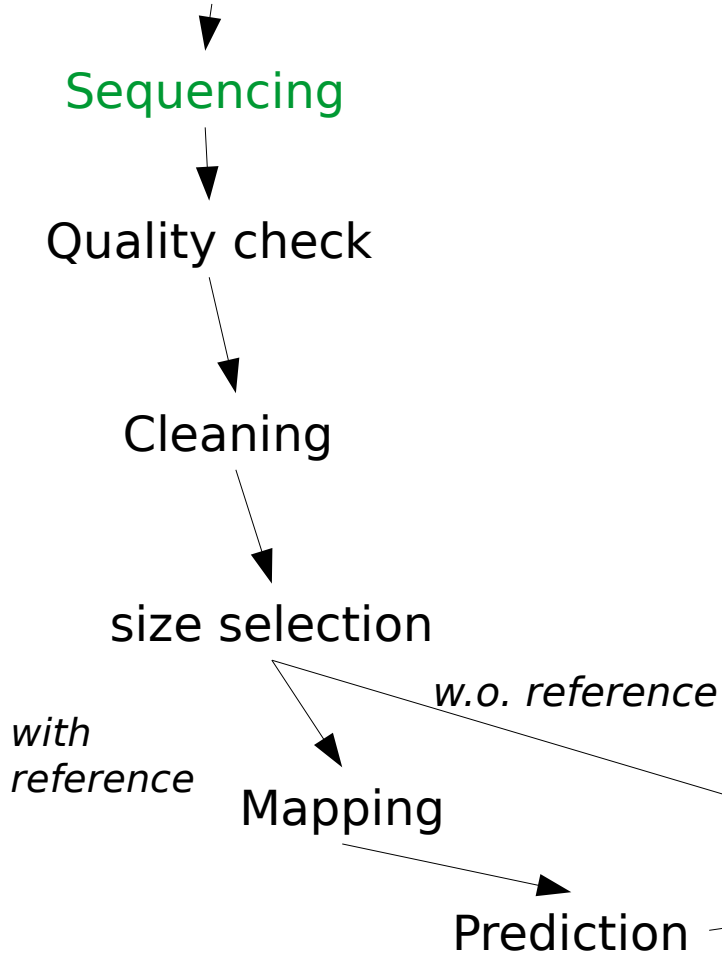*w.o. reference*

*with reference*

Mapping

Prediction

Annotation

Mapping: strict mapping

`Bowtie (not bowtie2), mismatch=1 max, considering multi best-hits`

Prediction: position in regard of annotation

`RSeQ, bedtools, s-mart, ...`

Module 3 : Normalization & differential analysis of RNAseq data

# smallRNAseq pipeline



Experimental design

Sequencing

Quality check

Cleaning

size selection

*w.o. reference*

*with reference*

Mapping

Prediction

Annotation

miRNA or not?

Module 3 : Normalization & differential analysis of RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

*w.o. reference*

*with reference*

Mapping

Prediction

Annotation

ex. of a classical miRNA

Module 3 : Normalization & differential analysis of RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

*with reference*

*w.o. reference*

Mapping

Prediction

Annotation

**Characteristics for 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

mature        star

Module 3 : Normalization & differential analysis of RNAseq data

# smallRNAseq pipeline

Experimental design

Sequencing

Quality check

Cleaning

size selection

*with reference*

Mapping

Prediction

*w.o. reference*

Annotation

Module 3 : Normalization & differential analysis of RNAseq data

**Software**:

**Database**:
- Rfam
- miRbase
- Silva
- GtRNAdb
- piRNA databank
...



miRanalyzer: a microRNA detection and analysis tool for next-generation sequencing experiments
Michael Hackenberg[1], Martin Sturm[1], David Langenberger[3,4], Juan Manuel Falcón-Pérez[6] and Ana M. Aransay[1,*]

BMC Bioinformatics

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]

Genome Biology

DSAP: deep-sequencing small RNA analys

miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades
Marc R. Friedländer[1], Sebastian D. Mackow

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

CPSS: a computational platform for the ana deep sequencing data

Discovering microRNAs from deep sequencing data using miRDeep
Marc R Friedländer[1], Wei Chen[1], Catherine Adamidi[1], Jonas Maaskola[1], Ralf Einspanier[3], Signe Knespel[1] & Nikolaus Rajewsky[1]

ep sequencing analysis

shortran: A pipeline for small RNA-seq data analysis
Vikas Gupta[1,2], Katharina Markmann[1], Christian N. S. Pedersen[2], Jens Stougaard a Andersen[1*]

miRDeepFinder: a miRNA analysis tool for deep sequencing of plant small RNAs
Fuliang Xie · Peng Xiao · Dongliang Chen · Lei Xu · Baohong Zhang

Hoeppner MP, Barquist LE, Gardner PP.: An introduction to RNA databases. Methods Mol Biol. 2014;1097:107-23

# smallRNAseq conclusion

Presented example: miRNA

Experimental design

Sequence characteristics to select the expected smallRNA

Annotation: first step of the biological analysis

RNAseq or high-throughput qPCR?

Akhtar MM, Micolucci L, Islam MS, Olivieri F, Procopio AD.: Bioinformatic tools for microRNA dissection. Nucleic Acids Res. 2016 Jan 8;44(1):24-44.

# RNAseq analyses

RNAseq: From sequence data (reads) to expression level (count)

Classical analyses of RNA-Seq:
- Check quality, Trimming
- Mapping / counts
- Assembly

Others usages of RNA-Seq:
- smallRNA study
- **expression at isoform level**

# Isoform

Eukaryotic gene includes introns:

exon

intron

donnor    acceptor
splice sites

gene locus

DNA

transcription + maturations
(splicing, coding gene +polyA tail, cap, etc)

mRNA                                              AAAA

One gene locus may rise diverse transcripts with different usages of exons

Alternative Splicing Event => isoforms

# Gene level

Leng N, Dawson JA, Thomson JA, Ruotti V, Rissman AI, Smits BM, Haag JD, Gould MN, Stewart RM, Kendziorski C.: EBSeq: an empirical Bayes hierarchical model for inference in RNA-seq experiments.Bioinformatics. 2013 Apr 15;29(8):1035-43

22

# Transcript level

Leng N, Dawson JA, Thomson JA, Ruotti V, Rissman AI, Smits BM, Haag JD, Gould MN, Stewart RM, Kendziorski C.: EBSeq: an empirical Bayes hierarchical model for inference in RNA-seq experiments.Bioinformatics. 2013 Apr 15;29(8):1035-43

23

# ASE-Alternative Splicing Events



Ding F, Cui P, Wang Z, Zhang S, Ali S, Xiong L. (2014) Genome-wide analysis of alternative splicing of pre-mRNA under salt stress in Arabidopsis. BMC Genomics, 15:431

24

# Organism specificity



AStalavista on the latests RefSeq versions
of species' annotations

Foissac S, Sammeth M (2007) ASTALAVISTA: dynamic and flexible analysis of alternative splicing events in custom gene datasets. Nucleic Acids Research 35:W297-299

- 5 common tools for DE
- RNAseq data:
    simulated data, *A.thaliana*
- number of common DE genes

Liu R, Loraine AE, Dickerson JA.Comparisons of computational methods for differential alternative splicing detection using RNA-seq in plant systems *BMC Bioinformatics*. 2014, Dec 16;15:364

26

# Comparing DE tools



- 5 common tools
- RNAseq data:
    simulated data, *A.thaliana*
- number of common DE genes

Hypotheses :
- level of analysis:
    transcript ≠ exon ≠ region
- organism specificity
- different methods/algorithms:
    mapping, counting, and DE

Tools for DE : >> 100, Tools for ASE : ~ 60, (2016)



ALEXA-Seq, Alt Event Finder, AltAnalyse, ARH-Seq, Asprofile, ASTALAVISTA, BitSeq, casper, Cufflinks/Cuffdiff, DEGSeq, DerFinder, DEXSeq, DiffSplice, DSGSeq, dSpliceType, ESFinder, eXpress, FDM, FineSplice, FlipFlop, FluxCapacitor, GliMMPS, GPSeq, iReckon, Iso-kTSP, IUTA, Jetta, JuncBase, KisSplice/KisDE, Limma, MATS/rMATS, Miso, MMSeq/MMDiff, PSGInfer, Quantas, rackJ, rDiff, Rmake, RNAprof, RSEM/EBSeq, rSeqDiff, SailFish, Salmon, SigFuge, Sircah, SNPlice, Solas, SplAdder, SpliceR, SpliceSeq, SpliceTrap, SplicingCompass, SplicingTypesAnno, SplicingViewer, SpliCQ, StringTie, Suppa, SwitchSeq

How to choose one ?          Benchmarking them !

# Benchmarking

Softwares?

- the number of tested methods is limited

Data?
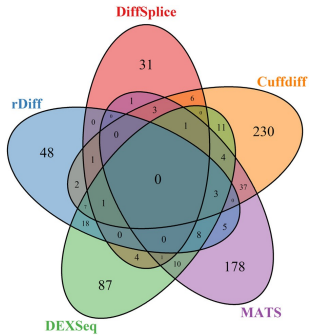
- RNAseq => how to know the truth?

    qPCR studies? Only small set of genes & cross-hybridization between isoforms

    RNA spikes? exogene sequences in controled quantity

- Simulated data => how to be as close as possible to the variability of the real data?
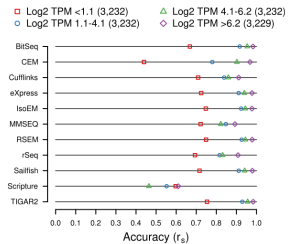
Task?

# Some benchmarking tasks
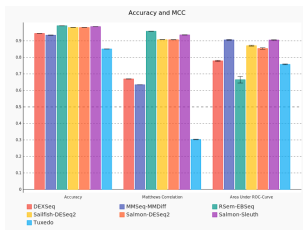


DE tools

How many common DE genes?



Isoform quantification

Expression rate, Exon number/transcript

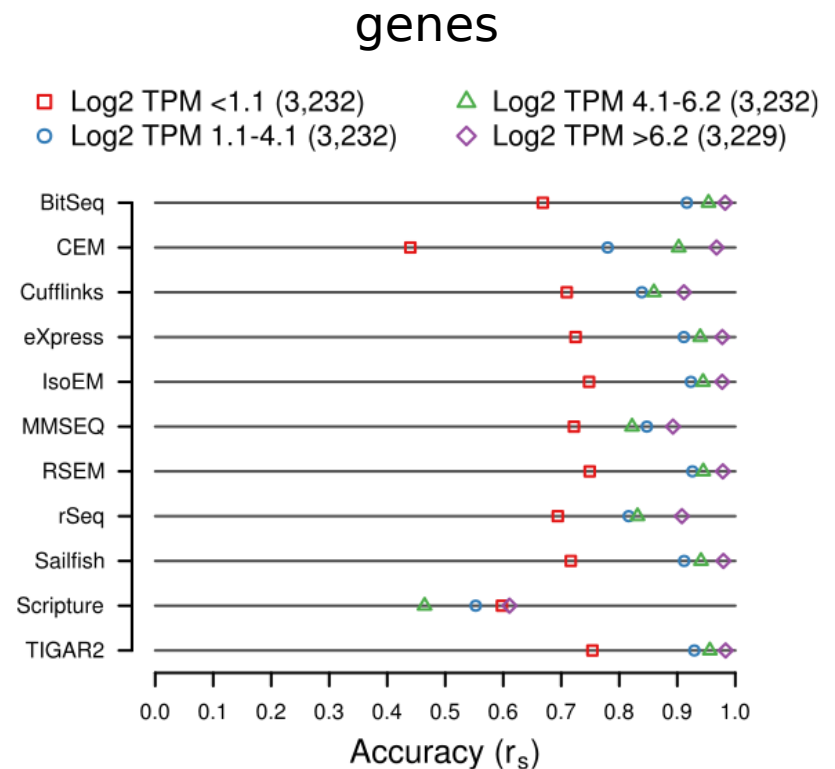Number of isoform/gene, ASE type



ASE detection

# Quantification: expression rate

Simulation:
Flux simulator
Human data set
RNAseq single-end
sequencing depth: 30 million reads

restricted on expressed transcripts
(10% of human transcripts)

Spearman correlation coefficient ($r_s$)
between the estimates and the
known input levels

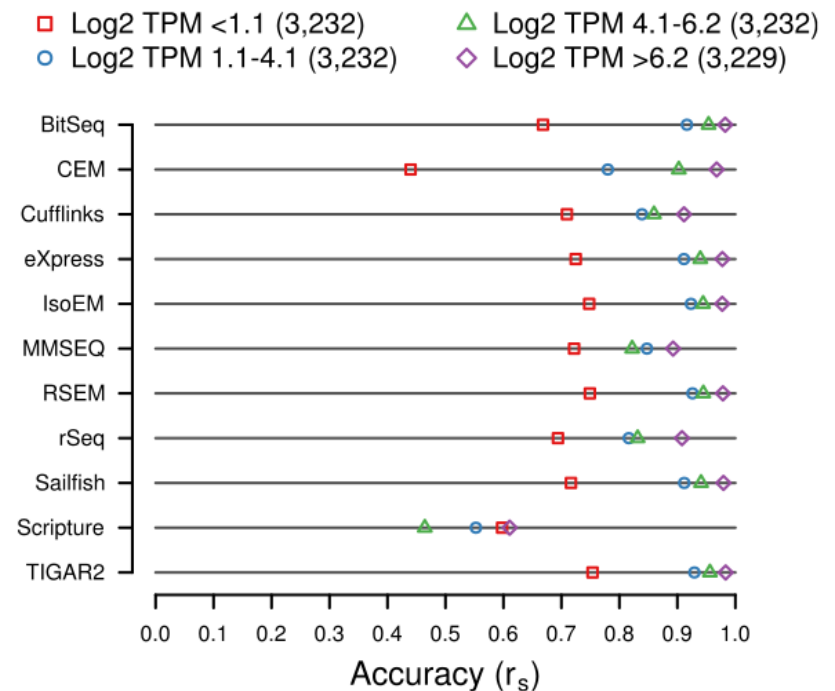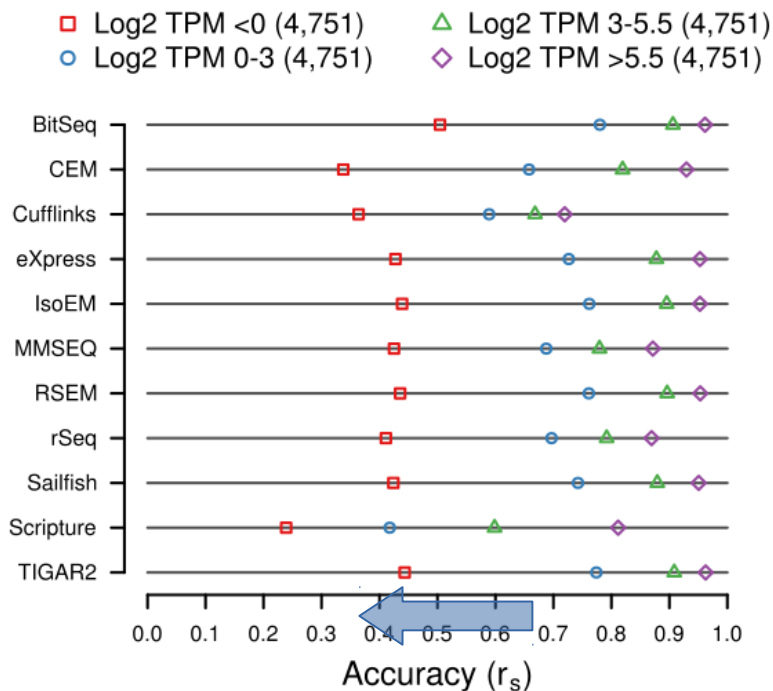4 bins of expression levels
(Log2 TPM)



genes

Kanitz A, Gypas F, Gruber AJ, Gruber AR, Martin G, Zavolan M. Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. Genome Biol. 2015 Jul 23;16:150

# Quantification: expression rate



transcripts

genes

Kanitz A, Gypas F, Gruber AJ, Gruber AR, Martin G, Zavolan M. Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. Genome Biol. 2015 Jul 23;16:150

32

# Quantification: exon number

#exons / transcript



Median expression levels:
0< Log2 TPM < 5.5

Cufflinks uses read-overlapping junction

Kanitz A, Gypas F, Gruber AJ, Gruber AR, Martin G, Zavolan M. Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. Genome Biol. 2015 Jul 23;16:150
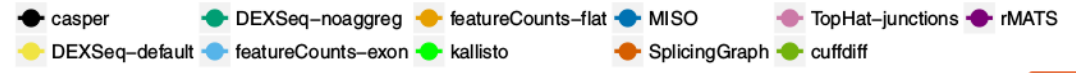
Simulation:
RSEM model from real data set + DTU for 1000 genes (switch of the relative abundances for the 2 most abundant isoforms between the conditions)

[i,j) i to j isoforms
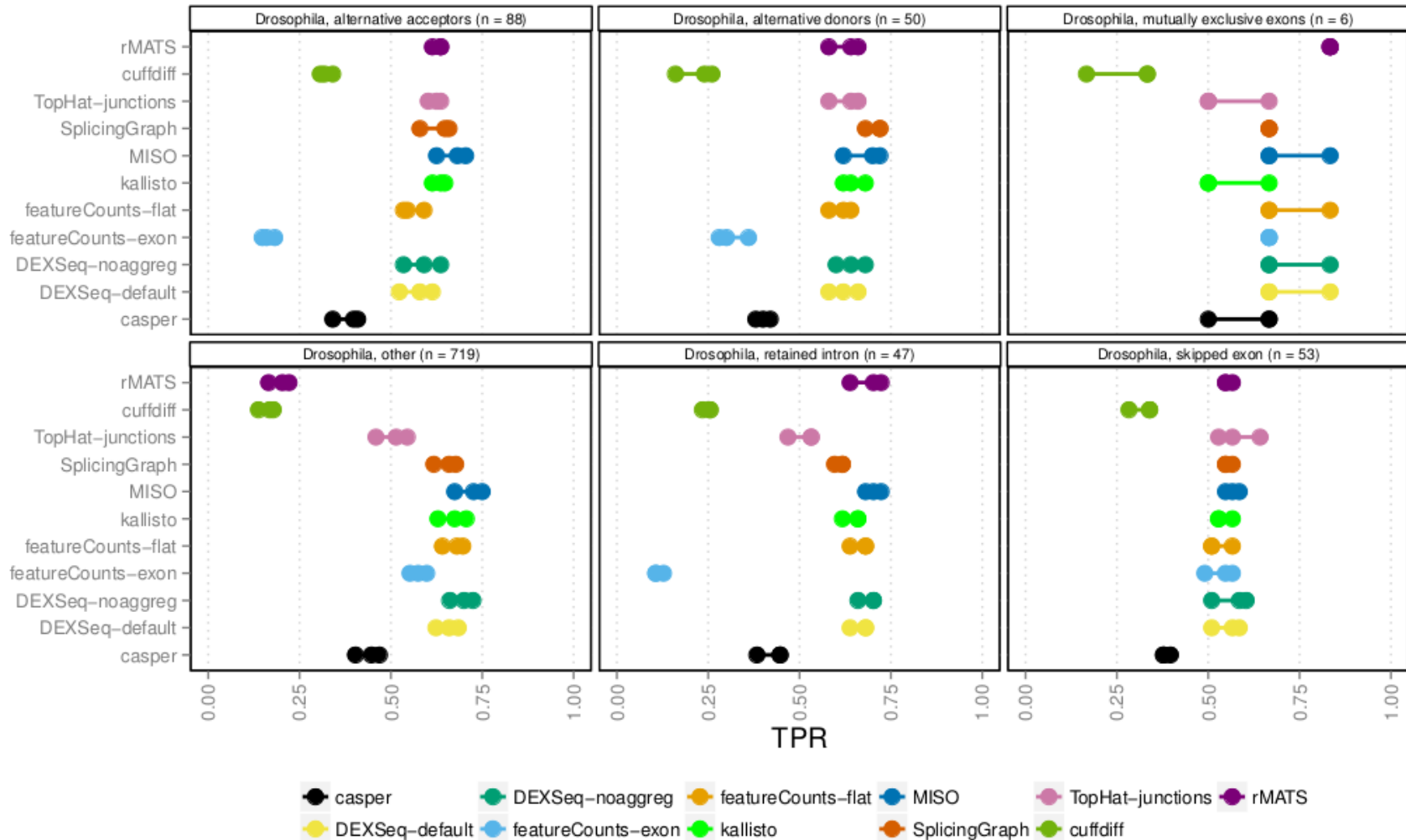n: gene number
n.ds: with DTU

3 circles, usual FDR thresholds (0.01, 0.05, 0.1): ideally, each circle should fall to the left of the corresponding vertical line

DTU: Differential Transcript Usage
FDR: False Discovery Rate
TPR: True Positive Rate



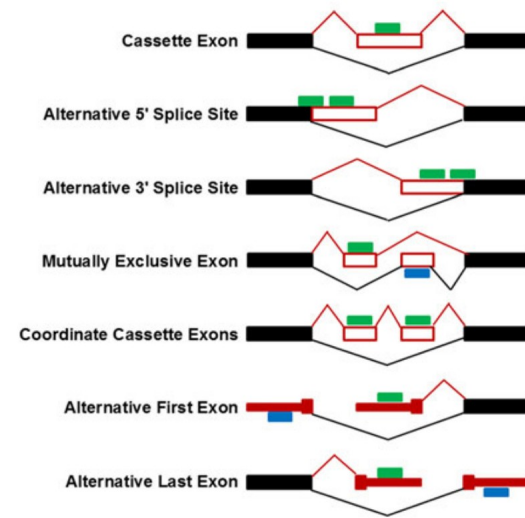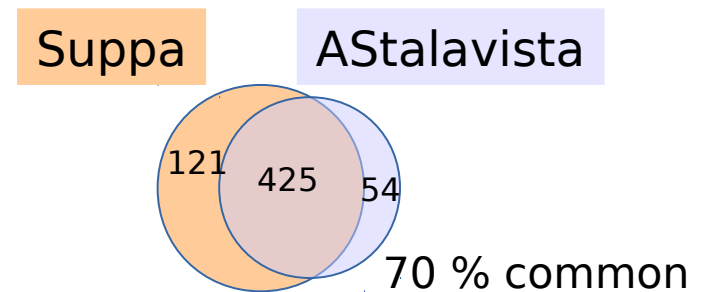Soneson C, Matthes KL, Nowicka M, Law CW, Robinson MD. Isoform prefiltering improves performance of count-based methods for analysis of differential transcript usage. *Genome Biol*. 2016 Jan 26;17(1):12

# Quantification: ASE ?

# Quantification: ASE ?

# ASE analysis

Annotation file
(gtf, gff, gbk)
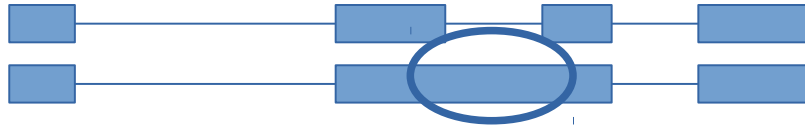gene
exon
transcript
CDS

→ ASE list

Suppa
AStalavista

Cassette Exon

Alternative 5' Splice Site

Alternative 3' Splice Site

Mutually Exclusive Exon

Coordinate Cassette Exons

Alternative First Exon

Alternative Last Exon

| chr22 hum. | Suppa | AStalavista |
|---|---|---|
| with ASE | 539 | 479 |
| without ASE | 733 | 793 |
| total | 1272 | 1272 |

Suppa    AStalavista

121    425    54

70 % common

Alamancos GP, Pagès A, Trincado JL, Bellora N, Eyras E. Leveraging transcript quantification for fast computation of alternative splicing profiles. *RNA*. 2015 Sep;21(9):1521-31
Foissac S, Sammeth M. ASTALAVISTA: dynamic and flexible analysis of alternative splicing events in custom gene datasets. *Nucleic Acids Res*. 2007 Jul;35(Web Server issue):W297-9
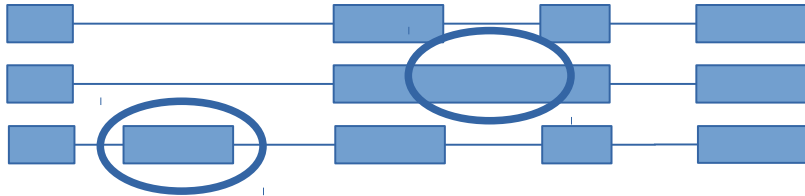
# ASE analysis need reference

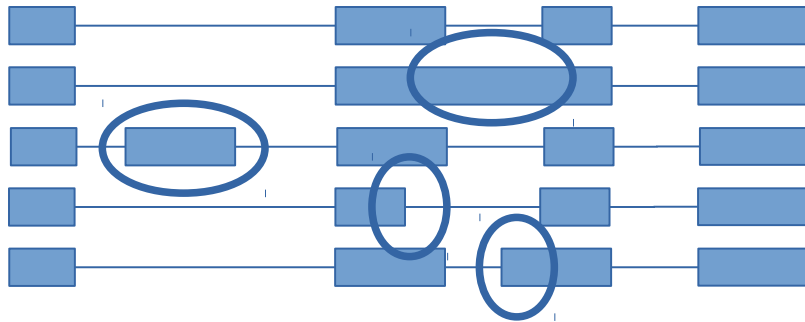# ASE analysis need reference

Retained Intron

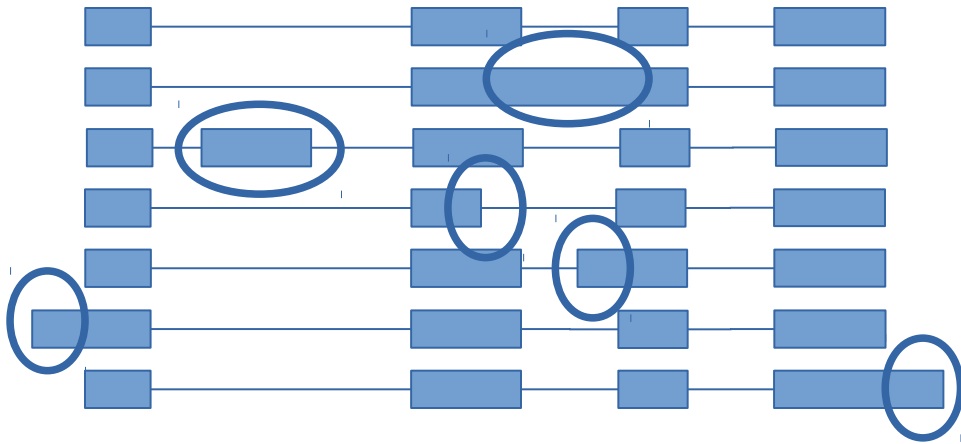# ASE analysis need reference



Retained Intron
Skipped Exon
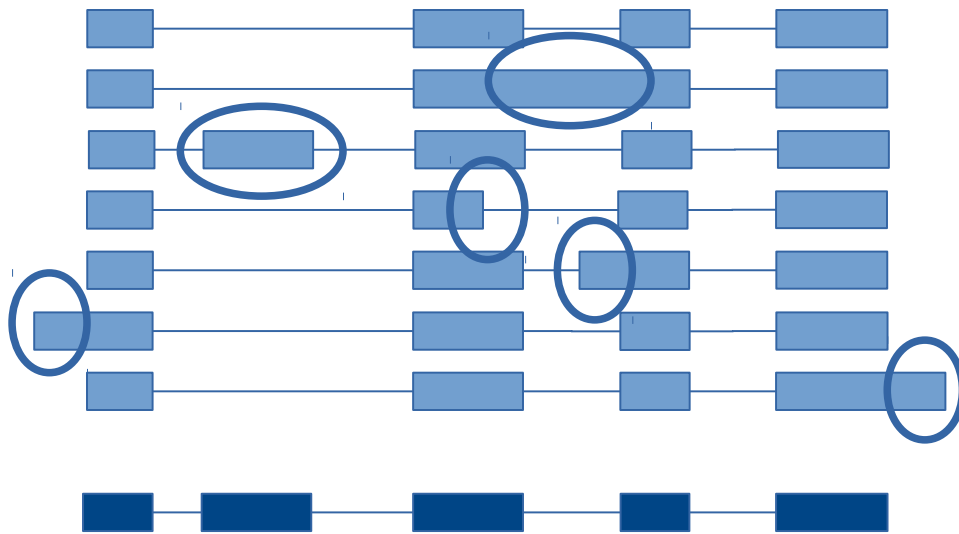
# ASE analysis need reference



Retained Intron
Skipped Exon
Alternative Donnor/Acceptor
Splicing Site

# ASE analysis need reference



Retained Intron
Skipped Exon
Alternative Donnor/Acceptor
    Splicing Site
Alternative First/Last Exon

# ASE analysis need reference



Retained Intron
Skipped Exon
Alternative Donnor/Acceptor
            Splicing Site
Alternative First/Last Exon

**Define a Reference Transcript** :
1 - **largest set** of **non-overlapping** exons
2 - that appends the **most frequently** among isoforms
3 - that covers the **widest area** over the gene region

=> may be a non "real" transcript
=> specific for each project

# Benchmark: ASE detection

Evaluate tools in their capacity to detect ASE from RNAseq data (neither the « right » rate of transcript expression, nor discovery of new expressed loci)

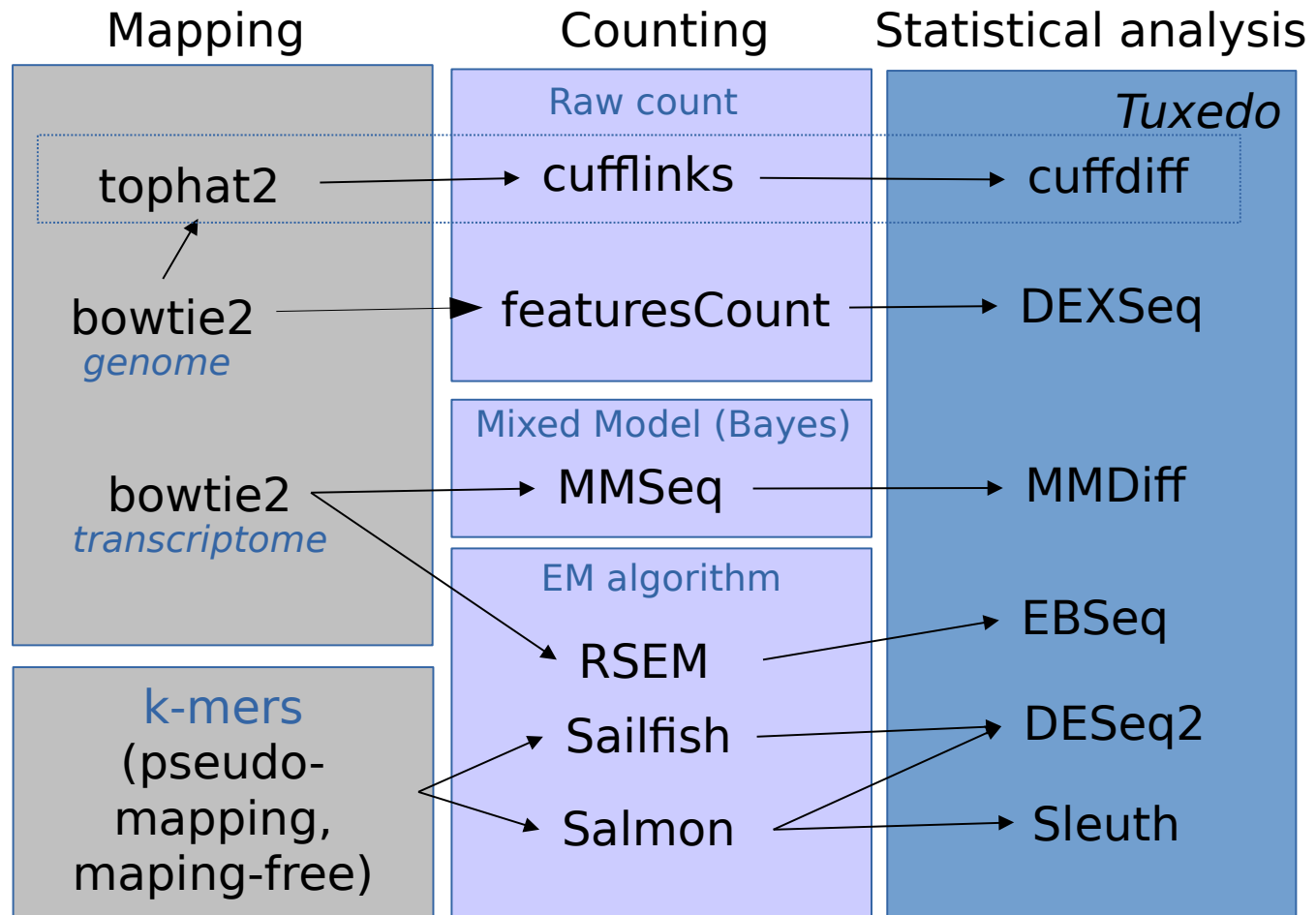Simulated data: controled expression rate of each isoform & the presence of each type of ASE :

- human chromosome 22, 744 genes
- 2 conditions x 3 replicates
- reads : pairend-end, 2 x 100 bp
- expression : 100 reads / transcript (no DE variation)
- For each type of ASE :
  - 10% of the transcript in one condition
  - Only reference transcripts in the other condition
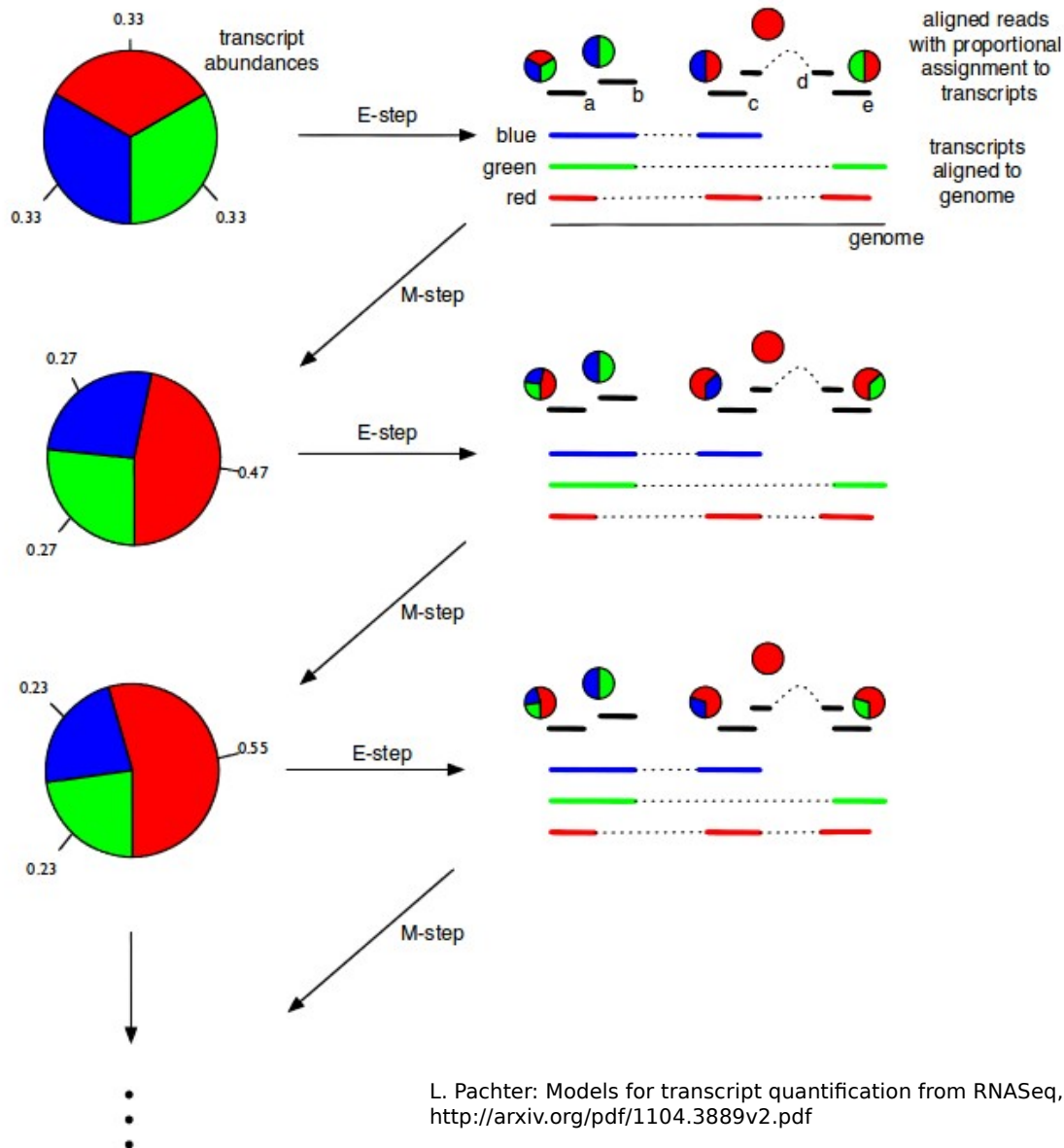  - 1 ASE/transcript/gene

Condition 1
100% references

Condition 2
90% references
10% with 1 ASE

Identification of a DE transcript => the method detects the ASE type

# Methods & tools benchmark

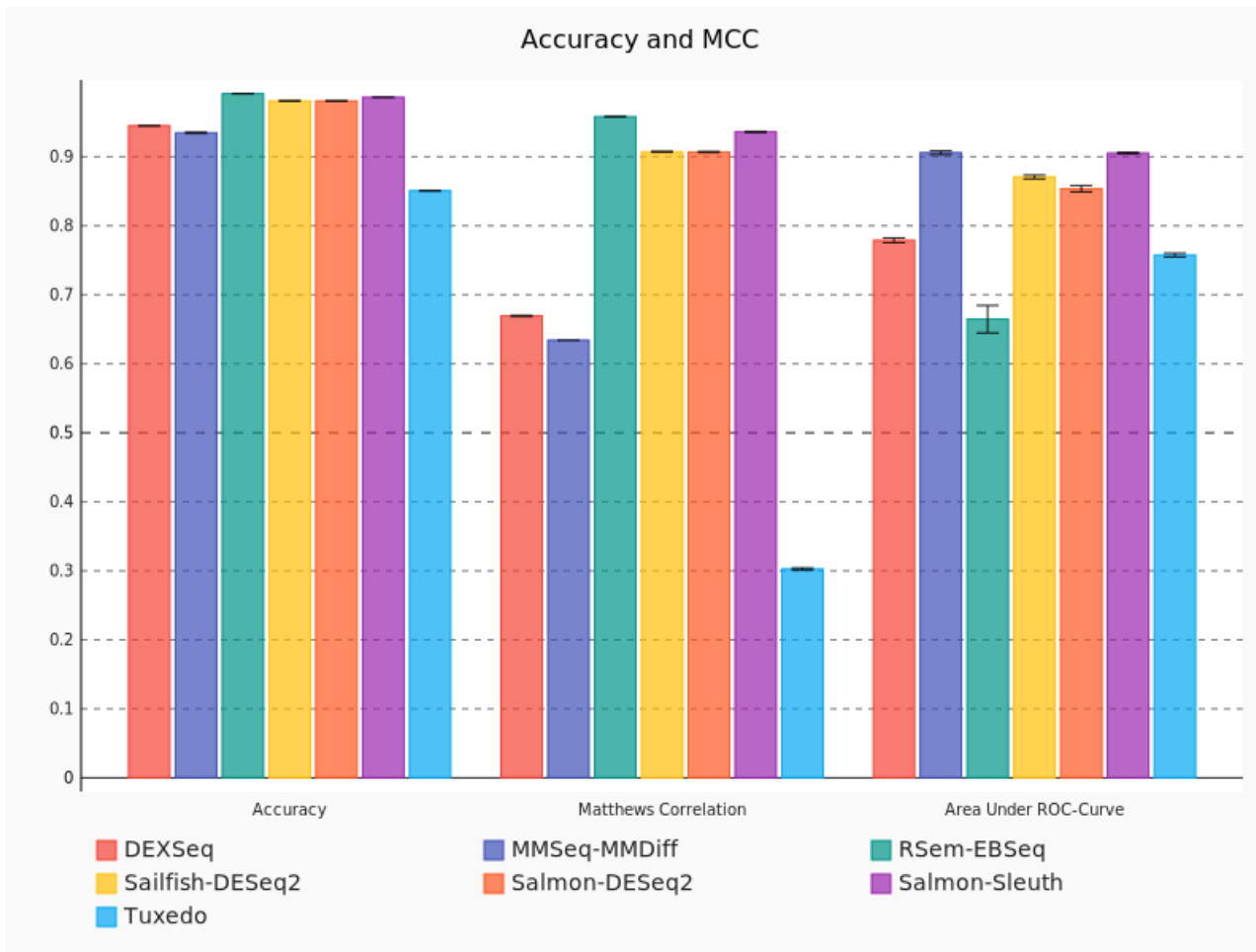Based-on the Expected-Maximisation algorithm

First 3 cycles of EM algorithm :

ex. : abundance of **red** transcipt estimated after the 1srt M-step:
(1/3 read a + 1/2 read b + 1 read d + 1/2 read e)/(total read number)

or (0.33+0.5+1+0.5)/5 =0.47

- proved to converge

- stop criterion implementation: when all probabilities that a fragment is derived from a transcript $\geq 10^7$ have a relative change of $\leq$ than $10^3$

L. Pachter: Models for transcript quantification from RNASeq, http://arxiv.org/pdf/1104.3889v2.pdf

46

# Results, alternative donnor site



Accuracy: are tool predictions correct?

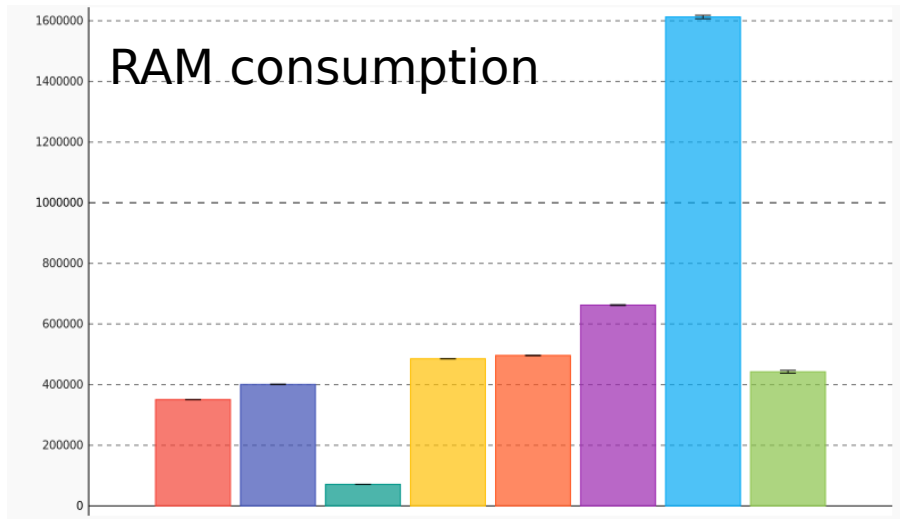RSEM/EBSeq, Sailfish, Salmon


MCC: if >0.5 then using tool is better than random
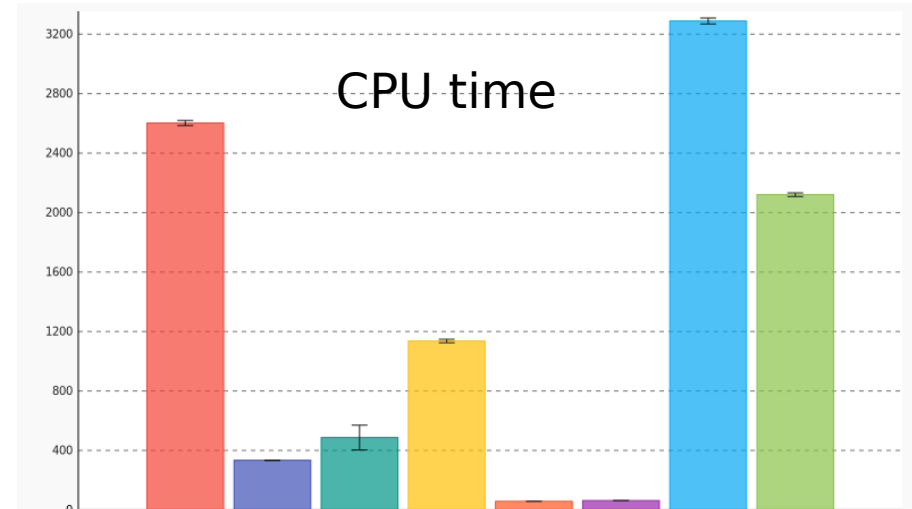
RSEM/EBSeq, Salmon-Sleuth
Tuxedo


AUC of ROC curve: confidence in the results of the tool

Salmon-Sleuth

# Performaces (ADss)

Max RAM usage (Ko)

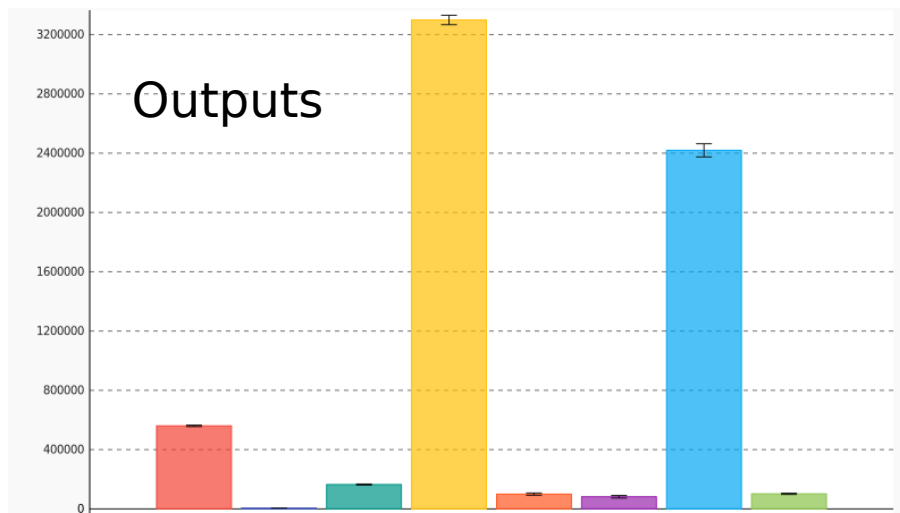RAM consumption

CPU time (seconds)

CPU time

Number of outputs

Outputs

- ■ Bowtie2
- ■ DEXseq
- ■ MMSeq/MMDiff
- ■ RSEM/EBseq
- ■ Sailfish/DESeq2
- ■ Salmon/Sleuth
- ■ TopHat2
- ■ Cufflinks/CuffDiff

Tuxedo : bowtie2+Tophat2+cufflinks/cuffdiff

# Choice of the method

For our benchmark

> human chr22, 10 % of ON/OFF transcripts with 1 ASE
> between 2 conditions, 3 replicates, 100 reads/genes

## Salmon/Sleuth

(RSEM/EBSeq, Sailfish)

developped pipeline: fastq ⇒ DE transcripts

appliance for IFB cloud, https://cloud.france-bioinformatique.fr

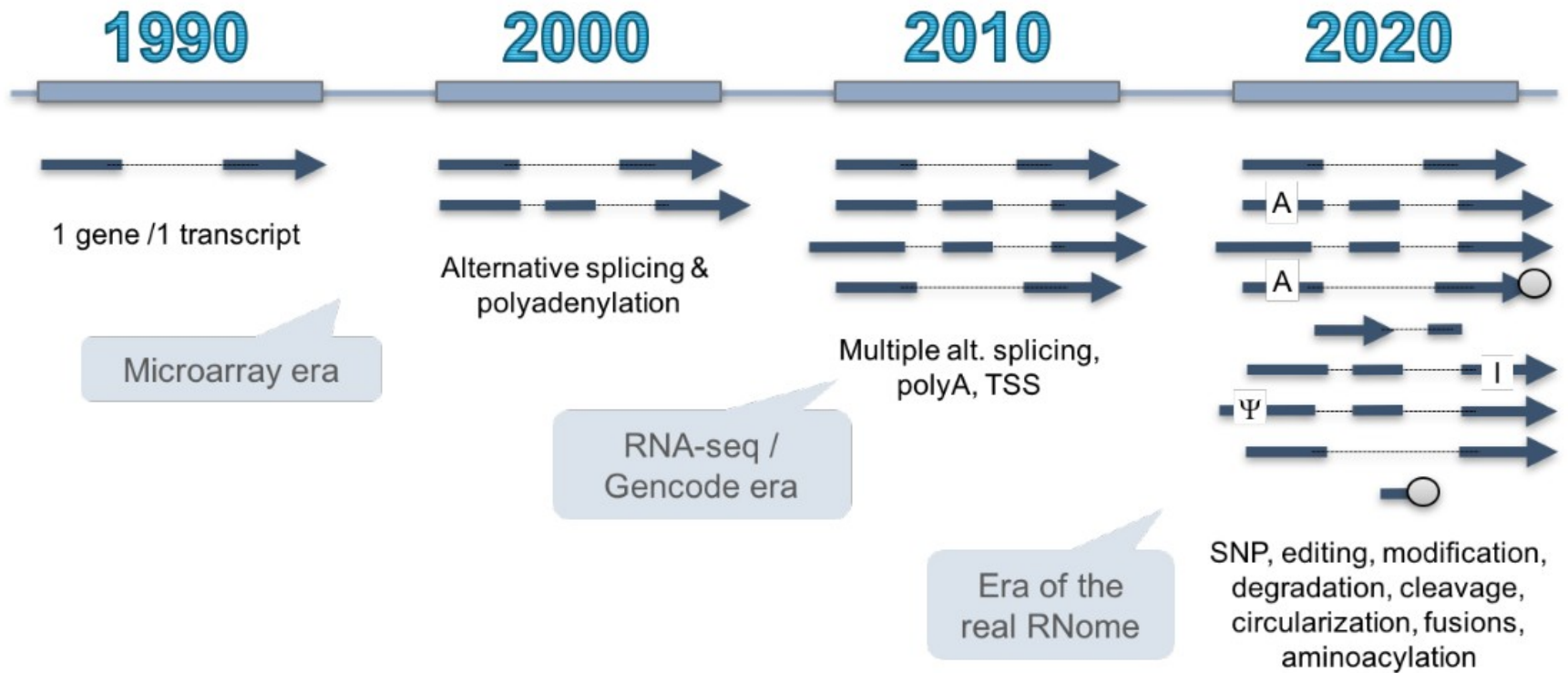# ASE Benchmarking conclusions

A benchmark is always limited:

➢ Software

➢ Simulation (human, plant, ...)
- Skipped exon, Retained Intron, Alternative First/Last Exon, Alternative 3'/5' Splicing Site
- 100 reads / gene (3-4 days, 6 cores, 13 G RAM)
  - 1000 reads / gene (to check if the wrong predictions result from the coverage deepness or from the algorithm)
- On/off condition : 1 ASE / transcript / gene :
  - more than 1 ASE / transcript
  - combination of ASE in the same transcript

Several benchmark studies should be conducted to have a global overview
The conclusions should be regularly update

# RNomics evolution

Part1-smallRNA: search for the method corresponding to the smallRNA
Part2-isoform level: will be easier with full length RNAseq technology



The more there is technical advances, the more we are going towards the unknown biology

# Thanks

**SPS — SACLAY PLANT SCIENCES**
Marie-Laure Martin-Magniette
Etienne Delannoy
Véronique Brunaud

**ifb    AVIESAN-IFB**

Eukaryotic small RNA
P. Bardou, C. Gaspin, S. Maman, J. Mariette, O. Rué, M. Zytnicki
http://www.france-bioinformatique.fr/sites/default/files/sRNA-Seq.pdf

Isoforms

Institut Pasteur    C3BI

**The SSFA team:**
Daniel Gautheret
Fabrice Leclerc
Jean Lehmann

**I2BC** — Institut de Biologie Intégrative de la Cellule

Bioinformatics and Biostatistics HUB
Marie-Agnès Dillies
Rachel Legendre
Hugo Varet

**eBIO**

Coline Billerey
Thibault Dayris
Marc Gabriel

UNIVERSITÉ PARIS SUD

université PARIS-SACLAY

cnrs
dépasser les frontières