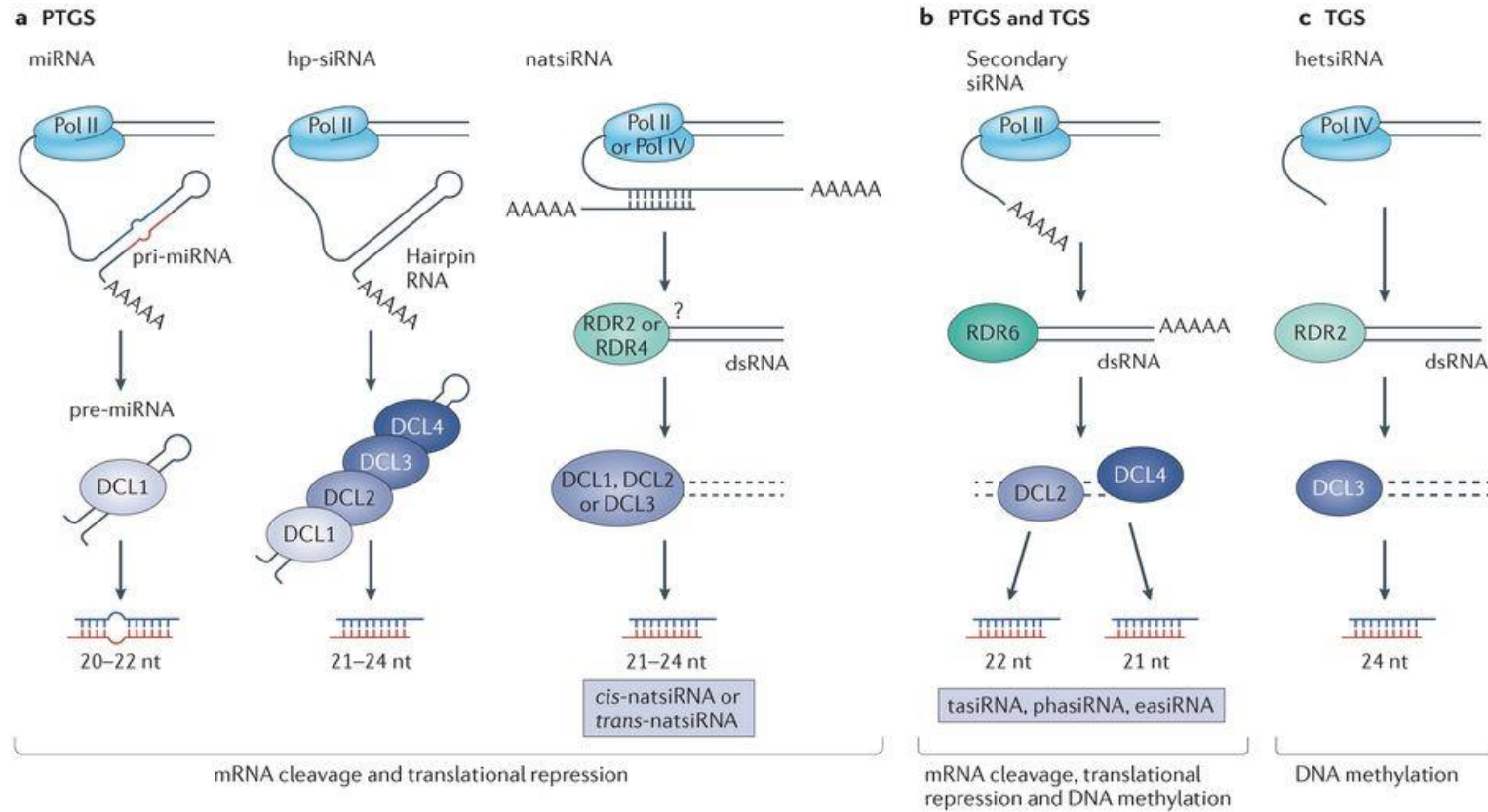


# Journée Annotation structurale et fonctionnelle des génomes eukaryotes : Annotation des smallRNAs

**Martine Da Rocha**

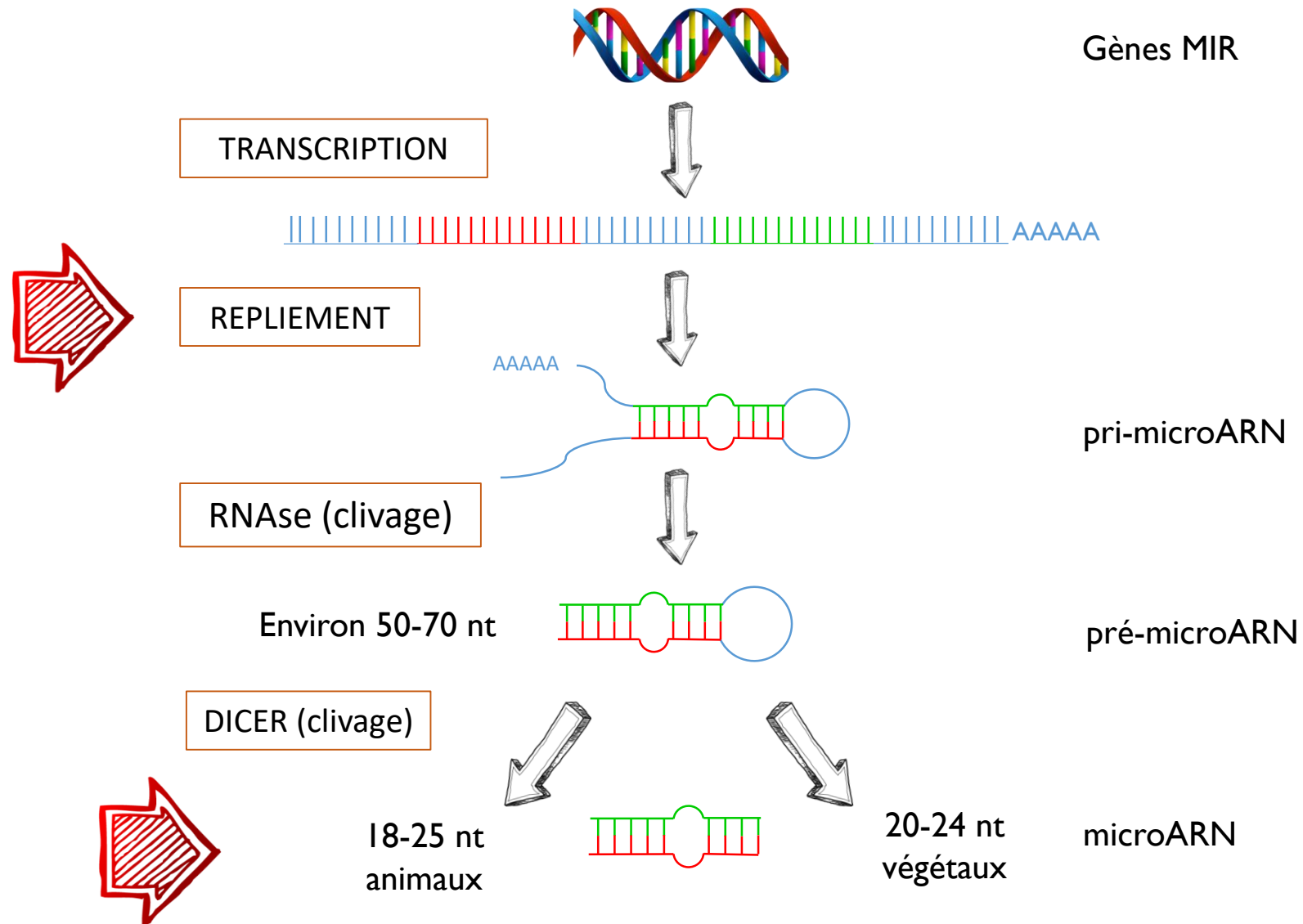
# SmallRNA Word



Nature Reviews | Molecular Cell Biology

Les SmallRnas -> Repliment -> maturation -> Structure action biologique

# Biosynthèse des microARN chez les eucaryotes



# Challenge de la prédiction des smallRNAs

Gènes codants pour des protéines

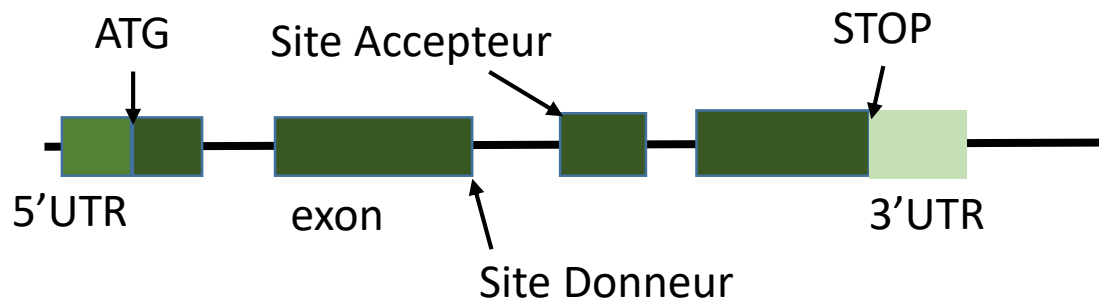


Schéma de structure lisible dans la séquence

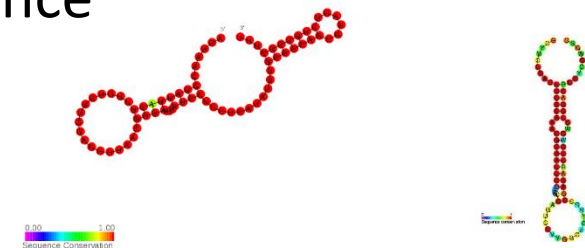
- cadres de lecture ouverts
- sites d'épissage
- fonction liée à la séquence

Gènes codants pour des SmallRNA



Absence de schéma de structure lisible dans la séquence

- fonction lié au repliement de la séquence



# Challenge de la prédiction des smallRNAs

- Critère de prédiction :
  - A partir de données de séquençage.
  - En fonction de la nature des petits ARNs : structure en hairpin
  - Pour les miRNAs à partir des miRNA connues (Conservation des séquences)

**problème beaucoup de faux positif**

Stratégies choisies :

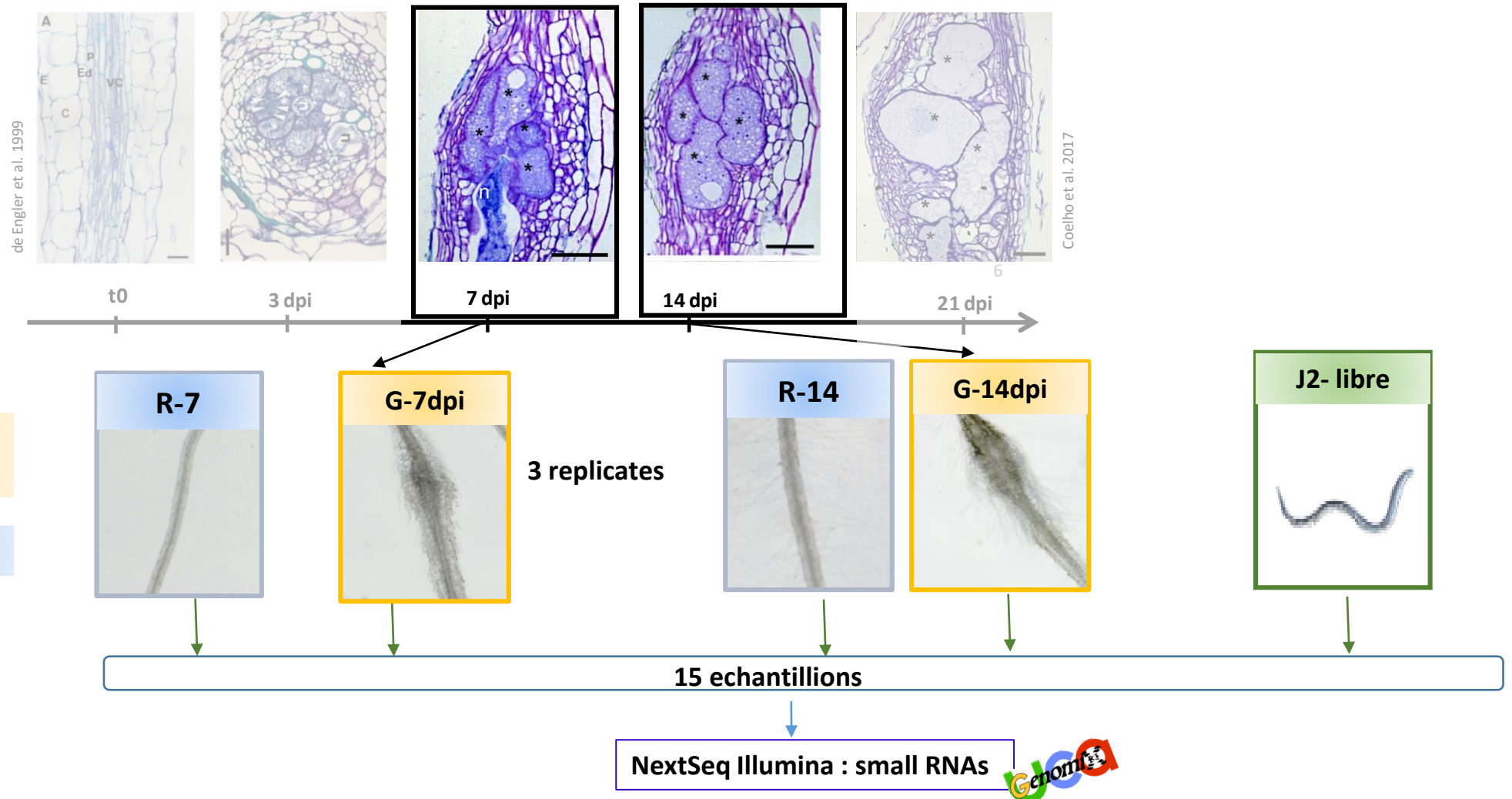
- Utilisation de répliquas biologique => même prédiction / dans les répliquas
- Si possible : recouper les résultats de plusieurs méthodes
- Attention au prédiction liée aux homologues : risque de propagés des erreurs

# Contexte de l'étude : Identification des sRNAs exprimés dans la racine de tomate lors de l'infection par un nématode

*M. incognita*



*Solanum lycopersicum*



de Engler et al. 1999

Coelho et al. 2017

Galles (G) : Racines infectées

Racines saine (R)

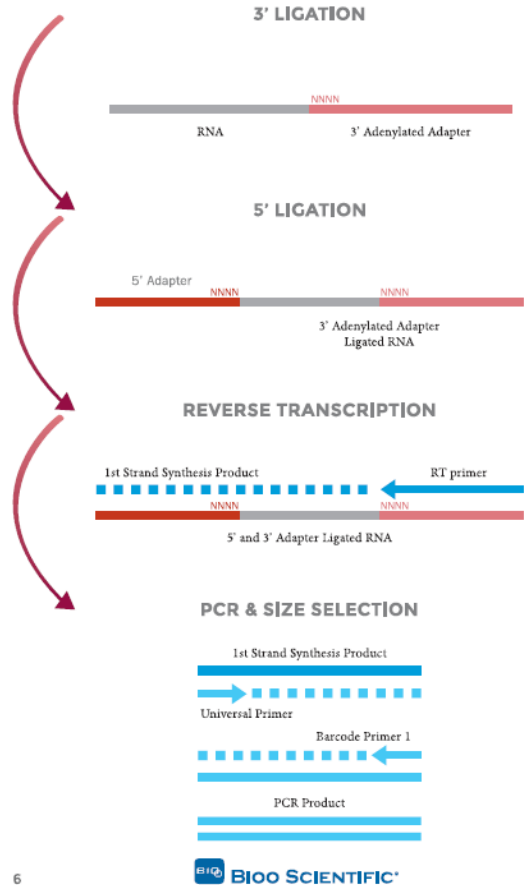
NextSeq Illumina : small RNAs



# Séquençage des smallRNAs

## NEXTflex™ Small RNA Sample Preparation Flow Chart

Figure 1: Sample flow chart.

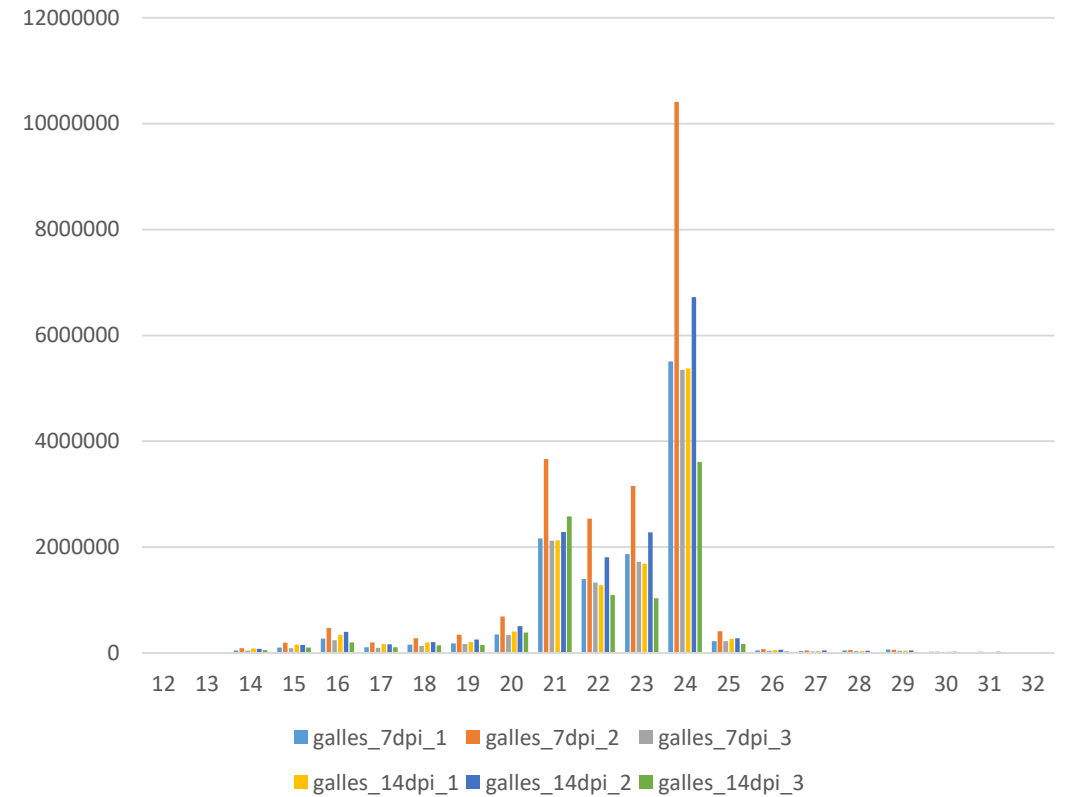
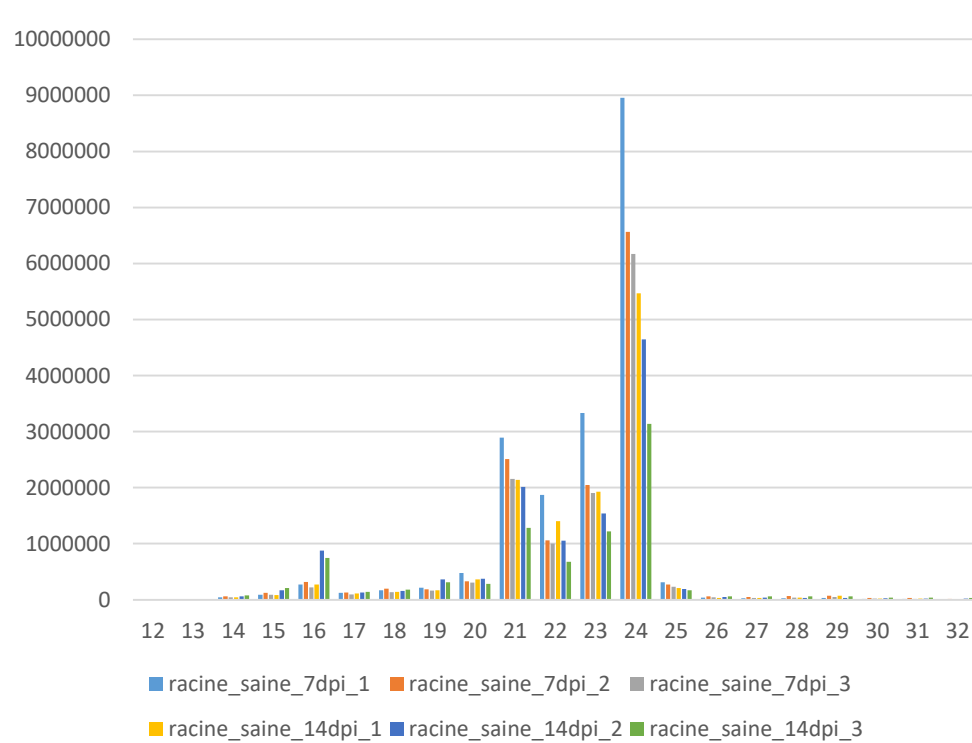


Purification basé sur la :

- Taille
- Adaptateur liaisons 5'P libres (ARNm ont une coiffe liaison impossible)

	Tissu végétal	Tissu animal	Nombre de reads
Référence Nématodes		J2_libres_1	38'251'985
		J2_libres_2	40'119'751
		J2_libres_3	34'527'534
Référence - Racines de 14 jours	Racines_saines_14dpi_1		13'971'658
	Racines_saines_14dpi_2		13'457'914
	Racines_saines_14dpi_3		10'334'113
Référence - Racines de 7 jours	Racines_saines_7dpi_1		20'995'892
	Racines_saines_7dpi_2		15'823'763
	Racines_saines_7dpi_2		14'167'136
Expérience - 14 jours après infection		Galles_14dpi_1	14'342'492
		Galles_14dpi_2	17'499'887
		Galles_14dpi_3	10'991'707
Expérience - 7 jours après infection		Galles_7dpi_1	14'083'472
		Galles_7dpi_2	25'164'984
		Galles_7dpi_3	13'419'506

# Tailles de fragments obtenues



Pics à 24 et 21 nucléotides



# Outils bioinformatique utilisés

- Tous les smallRNA : ShortStack

RNA. 2013 Jun;19(6):740-51. doi: 10.1261/rna.035279.112. Epub 2013 Apr 22.

**ShortStack: comprehensive annotation and quantification of small RNA genes.**

Axtell MJ<sup>1</sup>.

# Description de Shorstack

Script perl qui nécessite :

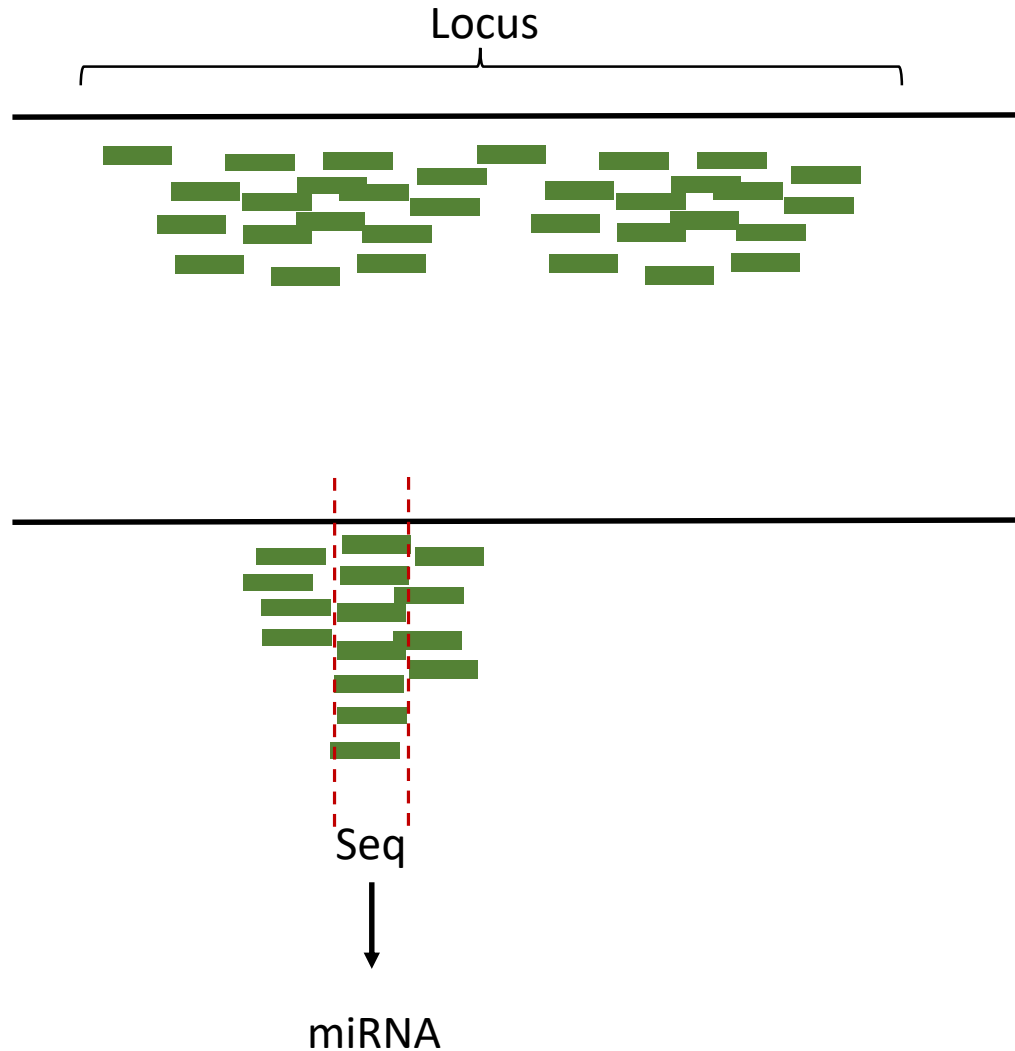
- Samtools
- RNALfold (Vienna RNA Package)
- RNAeval (Vienna RNA Package)

Utilisation Bowtie pour le mapping (peut être utilisé directement dans Shortstack) :

- Optimiser pour les reads <50 nt
- Supporte les données Solid.

Ce défini comme étant très stringent pour la détermination des miRNAs

# Fonctionnement de Shorstack



Identification de « ilot » -> Accumulation des reads

- Niveau de profondeur défini
- Extension de ces ilots aux ilots voisins (variation des smallRNAs le long des précurseurs)

## Création de Locus

- Comptage aux niveaux des locus de la zone la plus couverte : définition de la séquence la plus représentée Seq .
- Si  $20 < \text{Seq} < 24$  (plante) ou  $25$  (animal)  
Analyse par Rnafold -> Hairpin => miRNA

# Utilisation de Shorstack

1/Analyse de chaque réplique indépendamment.

2/Élimination de la redondance entre les banques / conservation seulement des locus présent dans 2 répliques

3/Relance Shortstack mais en lui donnant la liste des locus à déterminer



# Résultats de Shorstack

*Solanum lycopersicum*



Mapping :

~ 90% de lectures mappés

~ 30% mapper de façon multiple

locus prédit : 183632

miRNA : 111

*M. incognita*



Mapping :

~ 90% de lectures mappés

~ 70% mapper de façon multiple

locus prédit : 20210

miRNA : 91

Cohérents avec les résultats des analyses réalisées par Axtel, beaucoup plus de prédiction sur les plantes.

# Outils bioinformatique utilisés et dédiés au miRNA

Nucleic Acids Res. 2012 Jan;40(1):37-52. doi: 10.1093/nar/gkr688. Epub 2011 Sep 12.

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

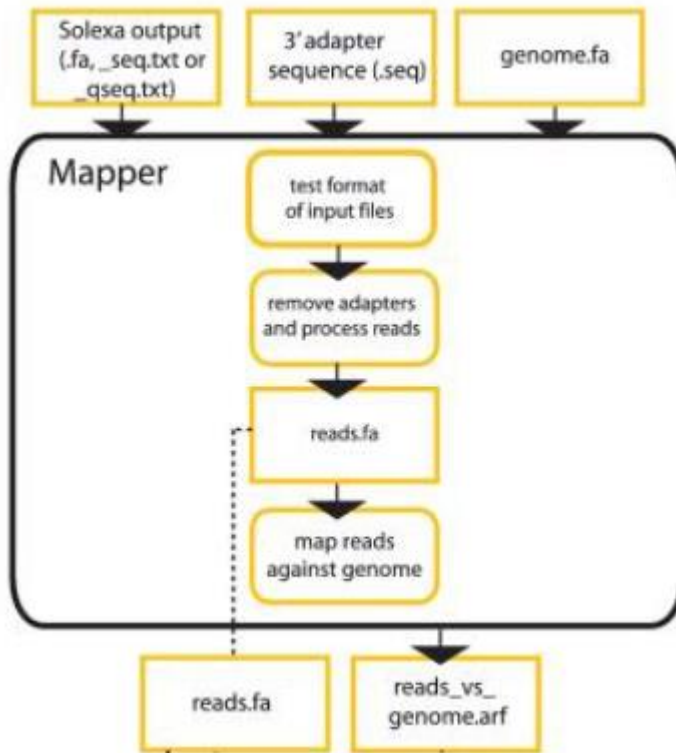
Friedländer MR<sup>1</sup>, Mackowiak SD, Li N, Chen W, Rajewsky N.

Bioinformatics. 2017 Aug 15;33(16):2446-2454. doi: 10.1093/bioinformatics/btx210.

**miRCat2: accurate prediction of plant and animal microRNAs from next-generation sequencing datasets.**

Paicu C<sup>1,2</sup>, Mohorianu I<sup>2,3</sup>, Stocks M<sup>2</sup>, Xu P<sup>3</sup>, Coince A<sup>3</sup>, Billmeier M<sup>3</sup>, Dalmay T<sup>3</sup>, Moulton V<sup>2</sup>, Moxon S<sup>3</sup>.

# Fonctionnement de Mirdeep2 : Module mapper



- Suppression des adaptateurs si nécessaire
- Les reads avec une séquence identique sont réduites pour supprimer la redondance. (Un nombre dans les nouvelles entêtes des fasta indique combien de fois la séquence est présente)
- Mapping avec par défaut 0 mismatch autorisé sur une graine de 18 nucléotides
- Fichier de sortie .arf et reads.fa



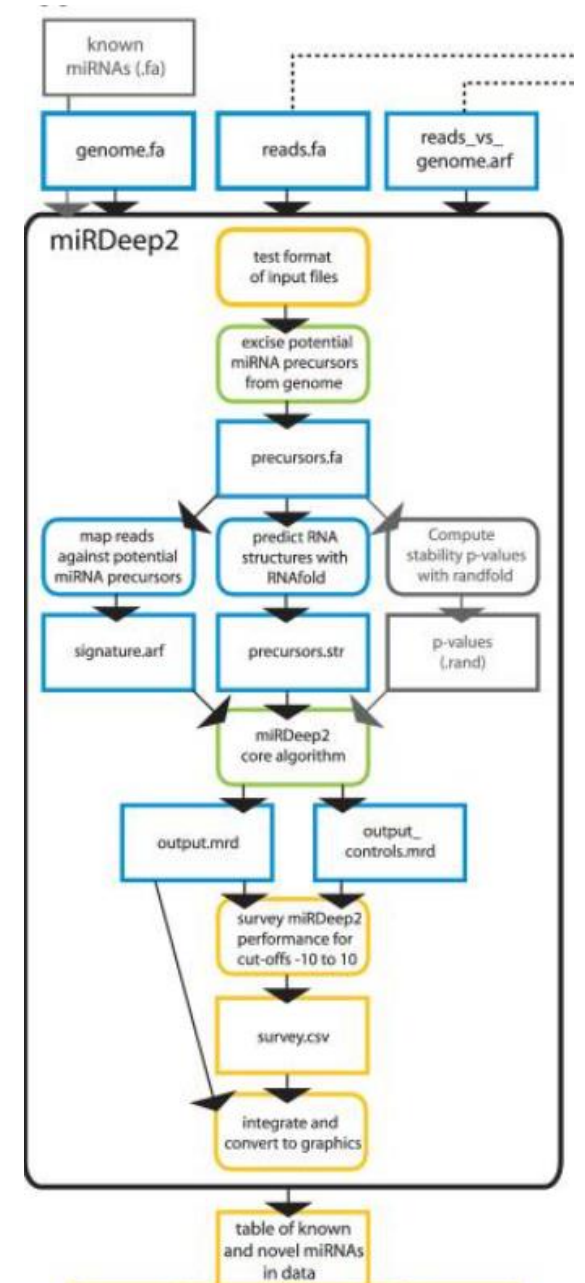
# Fonctionnement de Mirdeep2

Données d'entrées :

- Fichier de sortie .arf et reads.fa
- Génome
- miRNA connu

Etapas :

- Si les reads sont positionnés plus de 5 dans ce cas rejeté.
- Ne conserve que les alignements n'ayant aucun mismatch sur 18 nucléotides.
- Les zones contenant des reads mappés sont considérés comme précurseur potentiel et analysé avec RNAfold
- Recherche si les séquences de miRNA précurseurs connus se retrouvent sur le génome.
- Les miRNAs matures connus sont aussi positionnés sur le génome
- Statistique de performance de l'annotation est déterminé à partir des annotations déjà connus



# Adaptation mirDeep2 au plante

BMC Genomics. 2011 Feb 16;12:108. doi: 10.1186/1471-2164-12-108.

## Characterization of statistical features for plant microRNA prediction.

Thakur V<sup>1</sup>, Wanchana S, Xu M, Bruskiewich R, Quick WP, Mosig A, Zhu XG.

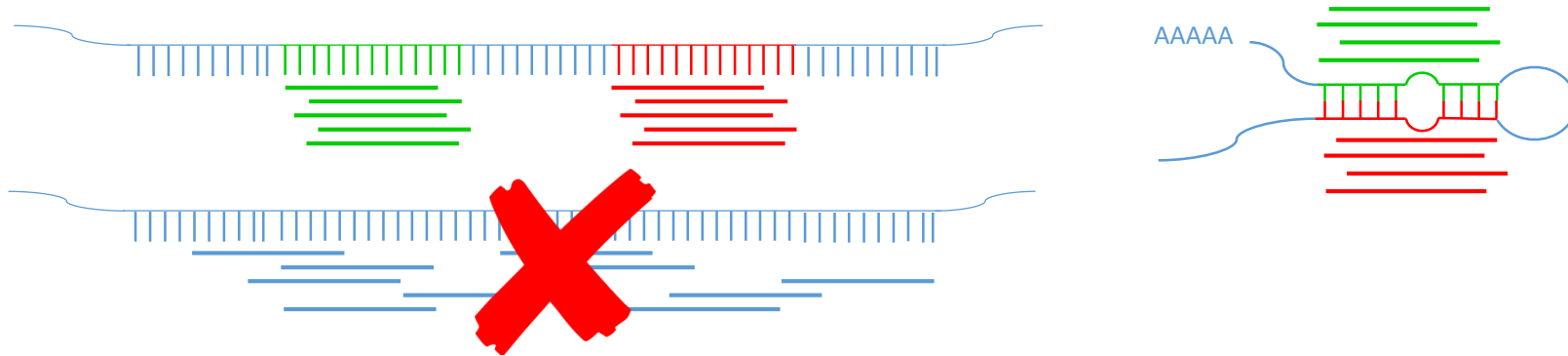
Feature		Original parameters of miRDeep	Plant specific (monocot)	Dicot specific (if any)
MFE	Known	Cumulative Distribution Function: $F(x) = e^{-e^{(x-\text{location})/\text{scale}}}$ Location = 32; Scale = 5.5	Log-odds score: $f(x)=a(b+e^{x^c})$	Log-odds score: $f(x)=a(b+e^{x^c})$
	Background	Cumulative Distribution Function: $F(x) = e^{-e^{(x-\text{location})/\text{scale}}}$ Location = 23; Scale = 4.8	a = 1.339e-12 b = 2.778e-13 c = 45.843	a = 4.46e-4 b = 9.125e-5 c = 26.929
Stability (log-odds)	Stable	1.6	1.37	0.63
	Unstable	-2.2	-3.624	-3.17
Nucleus conservation (log-odds)	Conserved	3	7.63	
	Non-conserved	-0.6	-1.17	
Excision length		140 nt	300 nt	
Paired	Total	≥14	≥15 nt	
Unpaired	Total	NA	≤5 nt	
	Consecutive unpaired	NA	≤3 nt	
Bulge	Total	could be as high as 5 nt #	≤2 nt	
Maximum multiple hits of deep-seq read		5	20	

**Table 1**  
List of miRDeep parameters estimated for plant specific miRNAs.

Modification du programme pour que celui corresponde au dicotylédone

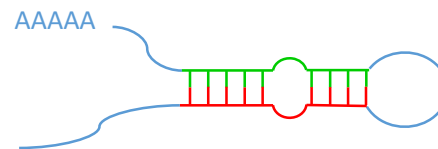
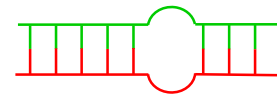
# miRCat2 : critères d'identification des microARN

- Sélection de zones candidates (accumulation d'alignements de lectures)



Utilisation de fenêtre coulissante

- Nombre maximal de positions différentes où s'alignent les lectures
- Taille des produits de clivage : *dicer call*
  - 18-25nt pour les animaux
  - 20-24nt pour les plantes
- Calcul de l'énergie de repliement



Variable en fonction du monde animal ou végétal

## miRCat2 :

Le UEA sRNA Workbench est disponible ici:

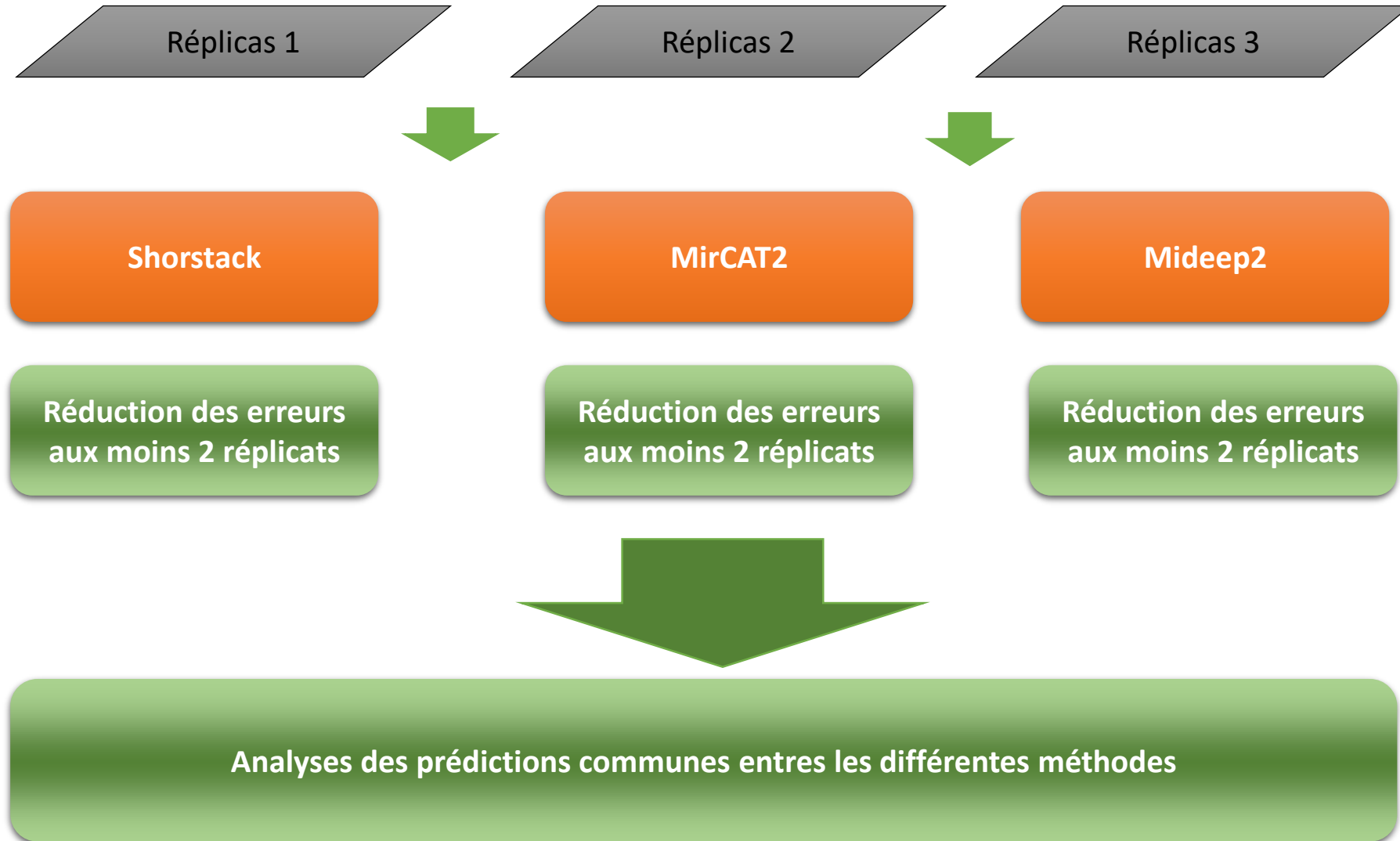
<http://srna-workbench.cmp.uea.ac.uk>.

Le code is available at:

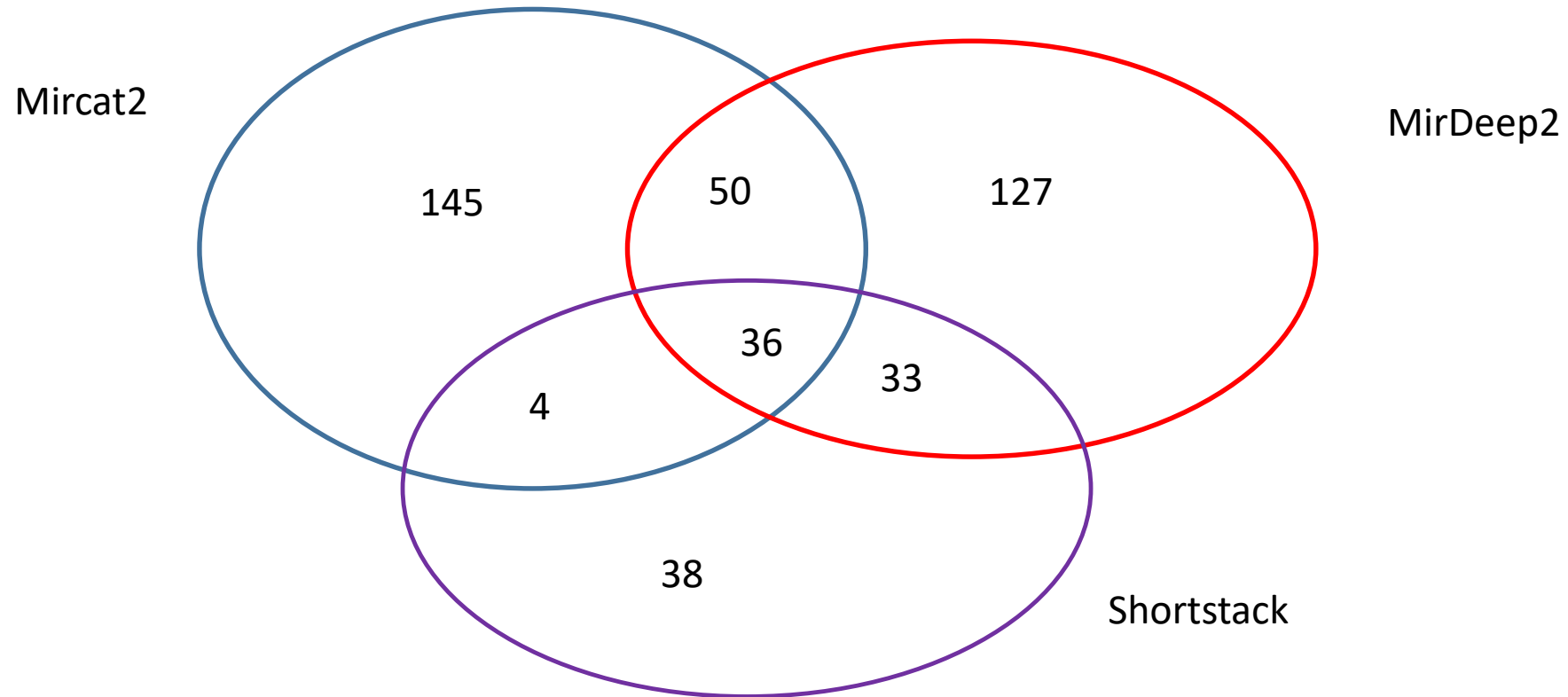
[https://github.com/sRNAworkbenchuea/UEA\\_sRNA\\_Workbench](https://github.com/sRNAworkbenchuea/UEA_sRNA_Workbench).

- Fichier de parametrage en fonction des caractéristiques des plantes /animals (format json)
- Galère d'utilisation via l'interface du Workbench.

# Pipeline d'analyse miRNA



# Résultats : Croisement des différents prédictions sur la tomate



Algorithme en cours de réécriture car particulier non pris en compte :  
Notamment inversion miRNA /Star miRNA

# Base de données sur les smallRNA

- [www.mirbase.org/](http://www.mirbase.org/) : Base de données dédiés au miRNA pour l'ensemble des espèces -> dernière mise à jour en Mars 2018 (après 4 ans sans mise à jours).
- DASHR – Database of small human noncoding RNAs ([www.lisanwanglab.org/DASHR/smdb.php](http://www.lisanwanglab.org/DASHR/smdb.php) );
- Base de données dédiées au small RNA chez le maïs et le riz (<http://sundarlab.ucdavis.edu/smrnas/>)
- Rfam base de données dédié au ARN non spécifique des smallRNAs



Stéphanie Jaubert-Possamai

Anais Vacquier



Virginie MAGNONE

Nicolas Pons

Kevin LEBRIGAND

