



# **Nouvelles approches optimisées pour le traitement de données sRNA-Seq**

Walid BEN SAOUD BENJERRI





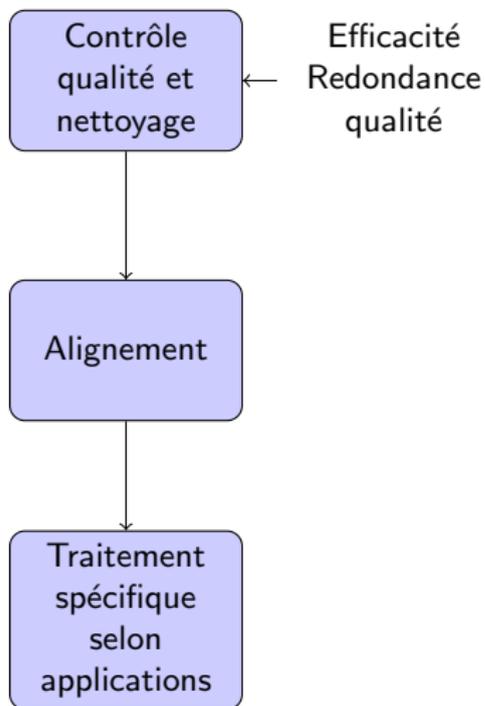
## Sommaire

- 1 Séquençage haut-débit et petits ARNs
- 2 srnacollapser
- 3 Conclusions et perspectives

## Génération de données

- Le génome = ensemble de séquences sur l'alphabet  $\{A,T,C,G\}$ , composé de plusieurs gènes, transcrits ou non.
- Le séquençage haut-débit
  - Assemblage de génomes
  - Analyse de variants (SNPs,...)
  - RNA-Seq : Accès aux gènes transcrits et leur niveau d'expression. En particulier on s'intéresse au sRNA-Seq : séquençage de petits ARNs.
- Verrou : Gros volume de données et redondance des lectures. Possibilité de synthétiser.

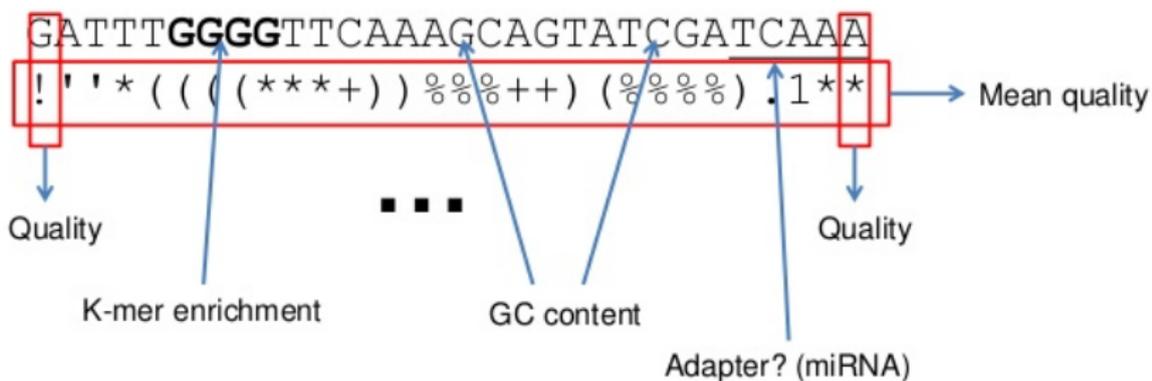
## Pipeline générique d'un traitement de données



## Exemple de fichier produit

```
Mock_A_TATAGCGA-GACACCGT_L001_R1_001.fastq
@HWI-M01141:63:A4NDL:1:1101:16668:1377 1:N:0:TATAGCGAGACACCGT
NACAGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCGCTAGGTGGTTTGT/
+
#>>AA>CAABBBGGGGGGGFFGHFEGGGFHHHGHFEGCEHHFEGGGGG@EEHHGGGHGHHK
@HWI-M01141:63:A4NDL:1:1101:14849:1418 1:N:0:TATAGCGAGACACCGT
NACGAAGGGTGCAAGCGTTACTCGGAATTACTGGGCGTAAAGCGTGCGTAGGTGGTGGTTT/
+
#>>>>A??AFAA1BGEGGAFFGCA0BFF1D2BCF/EEG/DBEE/E?GAEEFGAEFAEFG1
@HWI-M01141:63:A4NDL:1:1101:13802:1421 1:N:0:TATAGCGAGACACCGT
NACGGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCACGCAGGCGGTTTGT/
+
#>>AAABBBABBGGGGGGG?FGHGGGGHHHHHHHHGGGGHHHGGGGGGEGGGGGEGG?F/
@HWI-M01141:63:A4NDL:1:1101:15928:1426 1:N:0:TATAGCGAGACACCGT
NACGTAGGGTGCGAGCGTTAATCGGAATTACTGGGCGTAAAGCGTGCGCAGGCGGTTGTG/
+
#>>AABFB@FBBGGGGGGGGGGHGGGGFHHHHHHHHGGGGHHHGGGGGGGGGGGGGEE?G/
@HWI-M01141:63:A4NDL:1:1101:14861:1431 1:N:0:TATAGCGAGACACCGT
NACGAAGGGTGCAAGCGTTACTCGGAATTACTGGGCGTAAAGCGTGCGTAGGTGGTGGTTT/
+
#>>AAAABBFABGGGGGGEGEGHGGEFFHHHHHHHGHGGGGHHHGGGGGGEGFHHGGGGEGHE
@HWI-M01141:63:A4NDL:1:1101:15264:1465 1:N:0:TATAGCGAGACACCGT
NACGTAGGGTGCGAGCGTTGTCCGAATTACTGGGCGTAAAGAGCTCGTAGGTGGTTGTCC
+
```

## Une paire séquence/qualité



# MicroARNs

Les microARNs ont un profil d'expression abondant

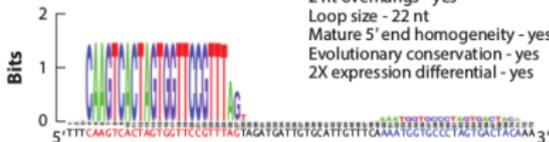
- Critères biologiques pour déterminer si une lecture est ou non un bras d'un mir

**a** hsa-mir-224 canonical mature and star

```

      CA          U U          A U  UG  UG
UUU  AGUCACUAG  GGU  CCGUUU  C  AGA  AU  \
AAA  UCAGUGAUC  CCG  GGUAAA  C  UUU  UA  U
      CA          -  U          A  -  GU  CG
    
```

Length of 5p arm - 24 nt  
 Length of 3p arm - 23 nt  
 Complementarity - 21 nt  
 2 nt overhangs - yes  
 Loop size - 22 nt  
 Mature 5' end homogeneity - yes  
 Evolutionary conservation - yes  
 2X expression differential - yes



**b** hsa-mir-212 canonical co-mature

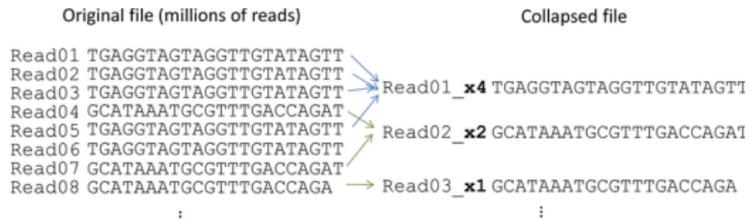
```

      -CA  U          CU          C          CCC  C
CGG  CC  UGGCU  AGACUG  UUAUCUG  GGG  \
GCC  GG  ACUGA  UCUGAC  AAUGAC  CCC  C
      ACC  C          CC          -          U--  G
    
```

Length of 5p arm - 23 nt  
 Length of 3p arm - 22 nt  
 Complementarity - 18 nt  
 2 nt overhangs - yes  
 Loop size - 17 nt  
 Mature 5' end homogeneity - yes  
 Evolutionary conservation - yes  
 2X expression differential - no



- Multiples séquences identiques, ces séquences mappent au même endroit, on économise du temps en mappant une seule séquence



- Mais les outils de mapping utilisent la qualité, on perd de l'information.
- Problème : Après avoir collapsé N séquences, que faire des N qualités :  $(Q_1, Q_2, \dots, Q_n) \Rightarrow Q_{synthetic} = ?$

- Hypothèse : N lectures identiques viennent d'un même locus, correspondant à une séquence (inconnue) notée  $s_1, \dots, s_n$ . On peut alors calculer une qualité à jour car pour une base donné une qualité supérieure discard toute qualité inférieure
- Ici la 1ère séquence nous dit que  $s_3$  est C avec 30% de probabilités, la 2ème dit la même chose avec 40% de probabilités, et la troisième également avec 50% de probabilités. Synthèse :  $P(s_3 = C) \geq 50\%$

ATCG	—30—
ATCG	—40—
ATCG	—50—

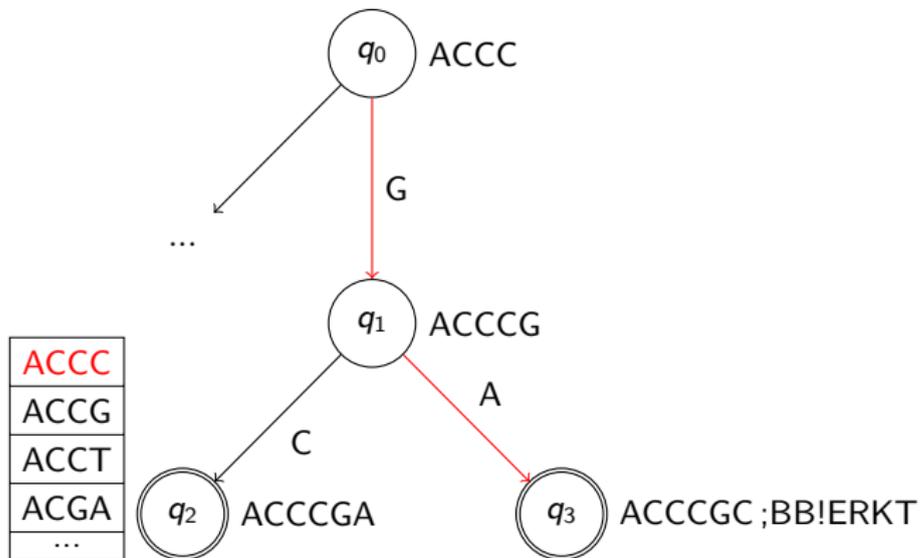
- Comptage des lectures par fichier.



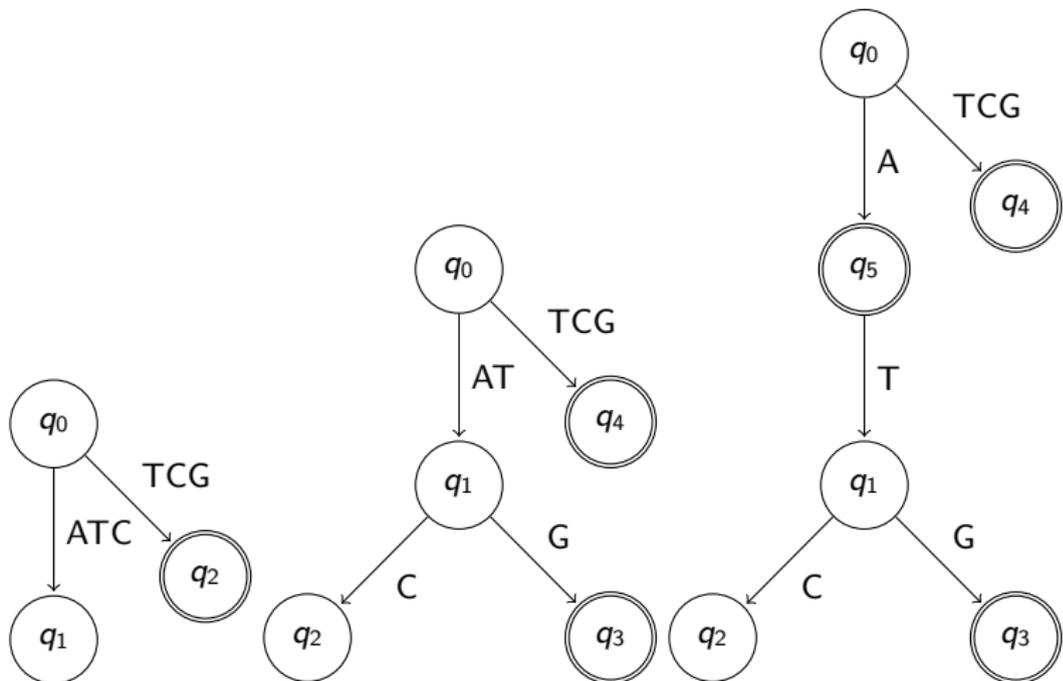
## Implémentation

- Méthode : utilisation de plusieurs tries (un pour chaque préfixe de séquence) : cela augmente la localité de la référence car les noeuds d'un trie sont stockés et lus ensemble.
- encodage binaire de l'ADN (2 bits overhead par nucléotide).
- les tries sont compacts, la recherche d'un LCP pour l'insertion se réduit à ctz (count trailing zeros) qui est une instruction directement implémentée dans les CPUs.
- implémenté en C++ pour un temps de calcul raisonnable même pour plusieurs gros fichiers. Optimizations bas niveau et multithreading pour optimiser l'usage de ressources.

## sRNACollapser tries



## Exemples d'insertions



## Résultats

	sRNACollapser	Naïve
A.Thaliana 1.7GB	16s	$\geq 2$ minutes
+ copy	+13s	$\geq 2$ minutes (and incrementing)



## Perspectives

- Correction d'erreurs
- Analyse différentielle
- Tableaux de hachage creux

## Perspectives : autres outils en cours de développement

- Mettre ensemble des mirs ayant supposément même fonction, comptabiliser erreurs de séquençage, séquences multimappés, comptage par famille au lieu de par locus pour éviter le rejet/biais
- Méthode : relier deux mirs si une certaine relation binaire est vraie