

Alignment, la spécificité d'abord !

Mikaël Salson

Équipe Bonsai, LIFL (Université Lille 1 et CNRS) – Inria

Les séquenceurs à haut débit

Les séquenceurs à haut débit

Roche 454

Illumina Hi-Seq

Ion Torrent

Solid

Les séquenceurs à haut débit

Roche 454



NIH

Illumina Hi-Seq

Ion Torrent

Solid

Les séquenceurs à haut débit

Roche 454

Haut débit
(\simeq 1 Gpb)

Ion Torrent

Illumina Hi-Seq

Très haut débit
($>$ 100 Gpb)

Solid

Les séquenceurs à haut débit

Roche 454

Moins de
une journée

Ion Torrent

Illumina Hi-Seq

Plusieurs jours

Solid

Les séquenceurs à haut débit

Roche 454

Illumina Hi-Seq

Lectures courtes
($\simeq 100$ pb)

Ion Torrent

Solid

Les séquenceurs à haut débit

Roche 454

Homopolymères

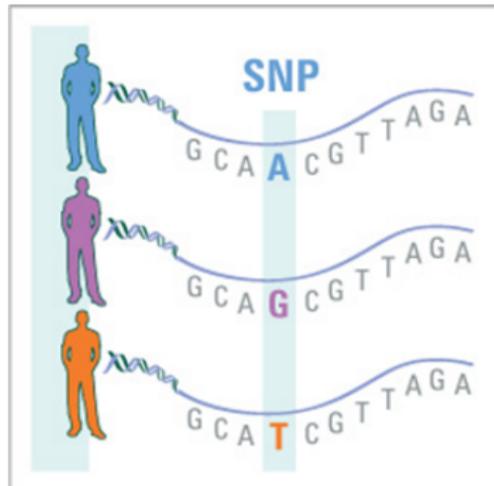
Ion Torrent

Illumina Hi-Seq

Substitutions

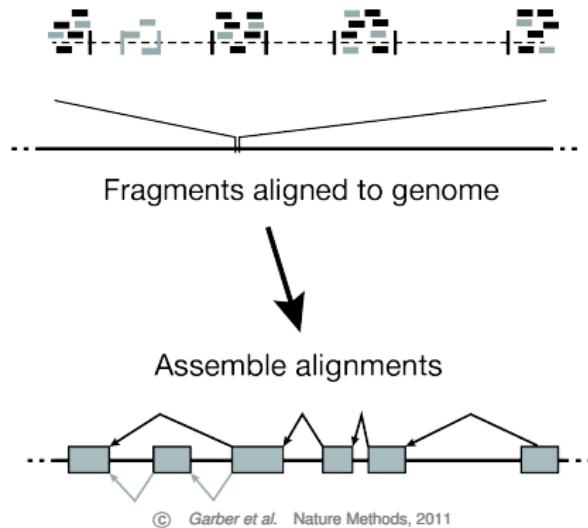
Solid

Le séquençage pour quoi faire ?

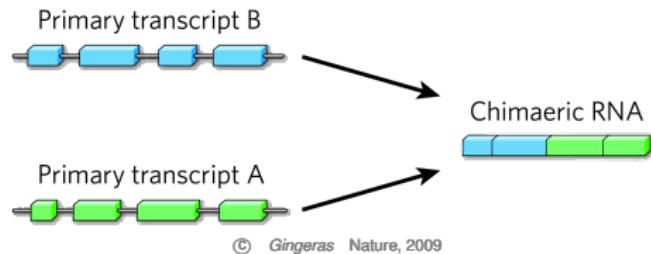


© Lauren Solomon Broad institute

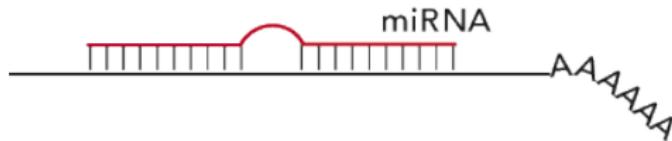
Le séquençage pour quoi faire ?



Le séquençage pour quoi faire ?



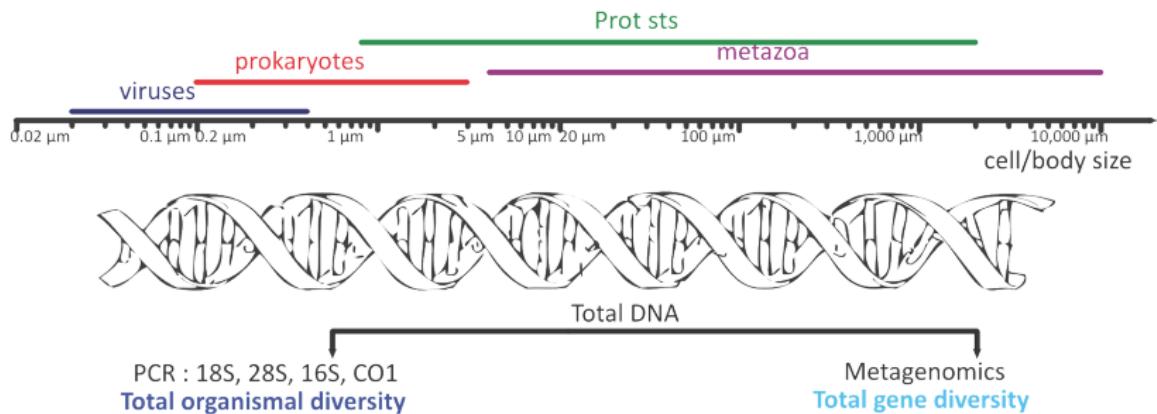
Le séquençage pour quoi faire ?



© Carin Cain, Victor Ambros Lasker foundation

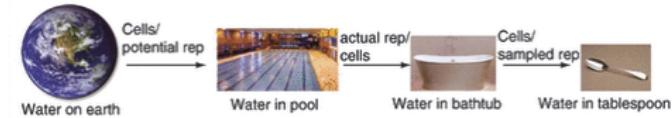
Le séquençage pour quoi faire ?

High Throughput Sequencing



CC BY Karsenti et al PLOS Biology, 2011

Le séquençage pour quoi faire ?



© Benichou et al Immunology, 2012

L'alignement de séquences

L'alignement de séquences

L'alignement de séquences

cagagtgactccatttgggtgag ?

L'alignement de séquences

cagagtgactccatttgggtgag ?

Recherche exacte

Knuth-Morris-Pratt

Boyer-Moore

Arbre des suffixes

Recherche exacte

Knuth-Morris-Pratt

Boyer-Moore

Années 1970

Arbre des suffixes

L'alignement de séquences

tgggcacctctatcttccg ?

L'alignement de séquences

tgggcacctctatcttccg ?

Recherche avec erreur(s)

Programmation dynamique

Shift-or

Recherche avec erreur(s)

Programmation dynamique

Années 1960–1980

Shift-or

Un problème compliqué ?

Knuth-Morris-Pratt

Boyer-Moore

Shift-or

Programmation
dynamique

Arbre des suffixes

Un problème compliqué ?

Knuth-Morris-Pratt

Boyer-Moore

Trop de temps

Shift-or

Programmation
dynamique

Arbre des suffixes

Un problème compliqué ?

Knuth-Morris-Pratt

Boyer-Moore

Shift-or

Trop de temps

Programmation
dynamique

Arbre des suffixes

Trop d'espace

Un problème compliqué ?

Beaucoup de données

Spécificités du protocole et du séquençage

Graines

Génomie

Graines

Génome



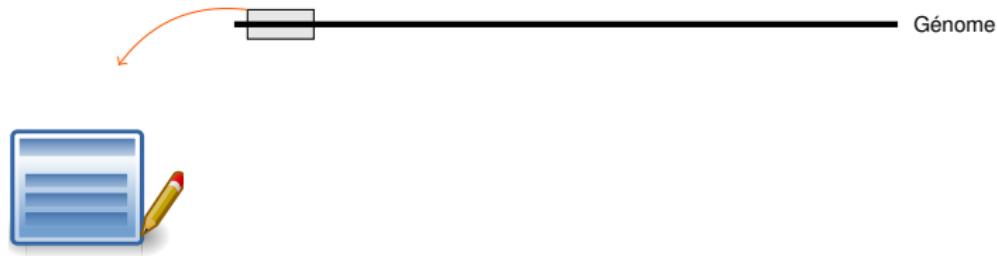
CC BY SA RRZEicons Wikimedia

Graines



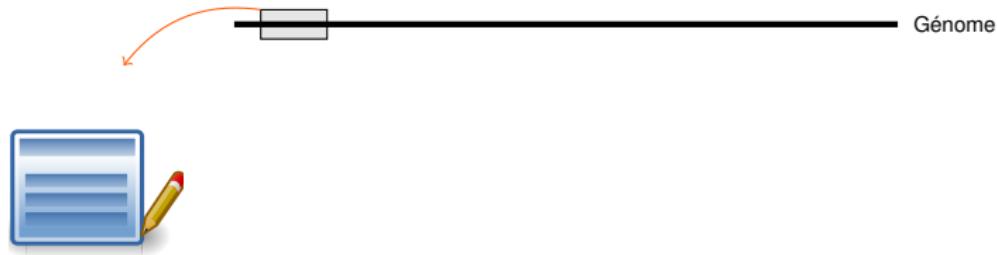
CC BY SA RRZEicons Wikimedia

Graines



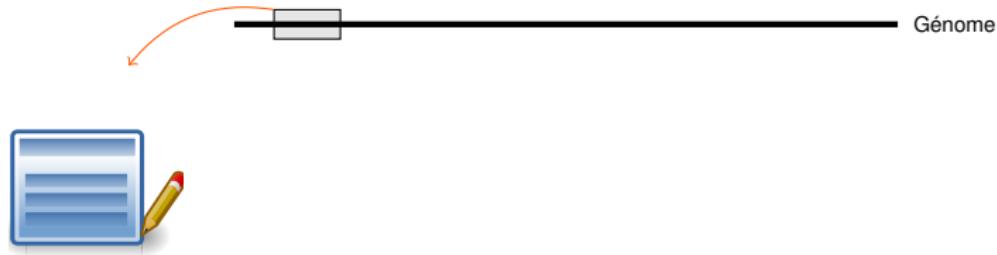
CC BY SA RRZEicons Wikimedia

Graines



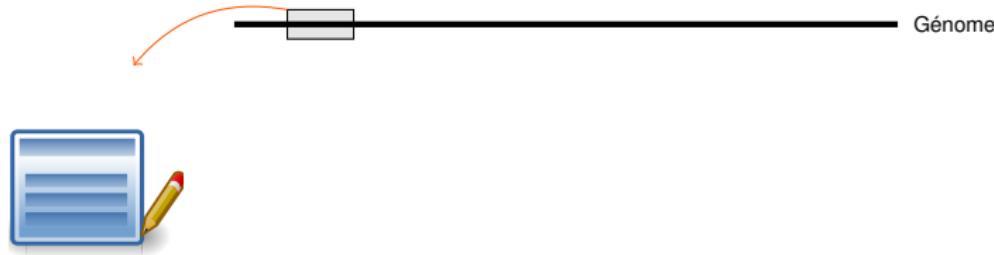
CC BY SA RRZEicons Wikimedia

Graines



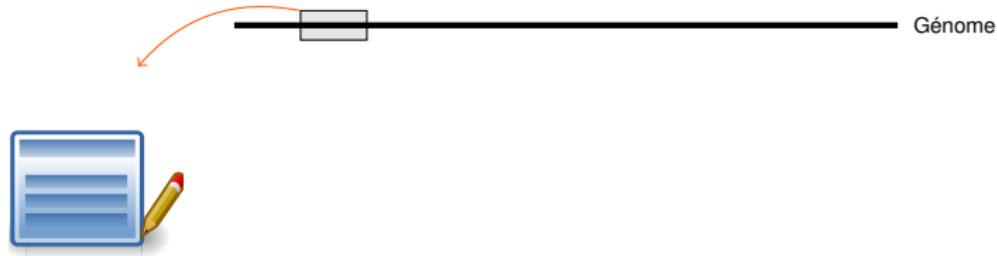
CC BY SA RRZEicons Wikimedia

Graines



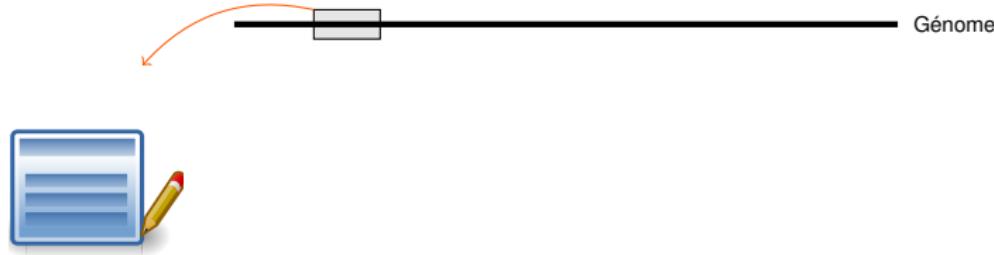
CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



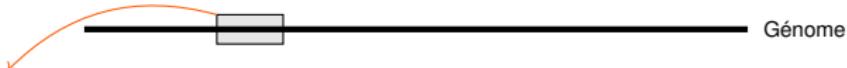
CC BY SA RRZEicons Wikimedia

Graines



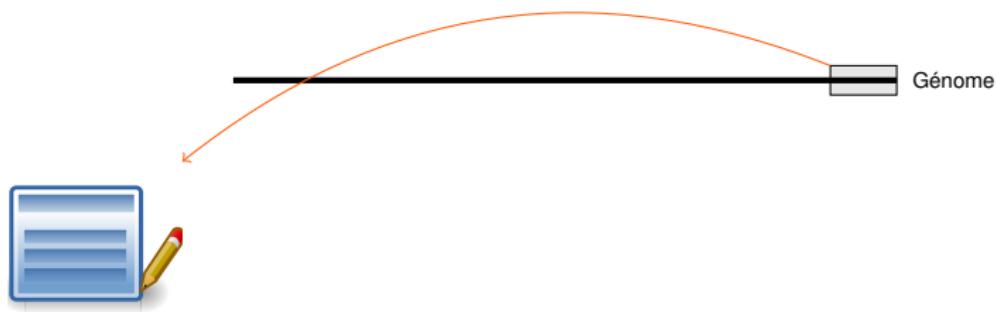
CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines

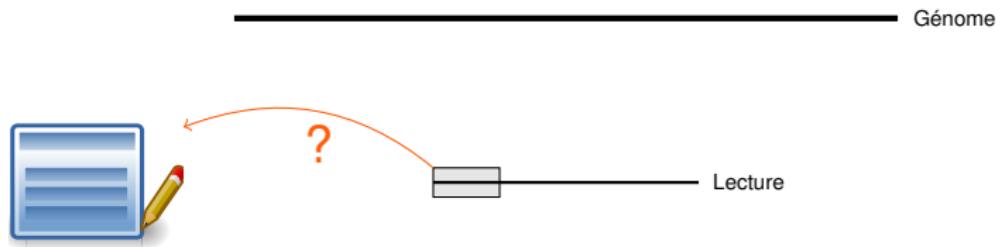
Génome



Lecture

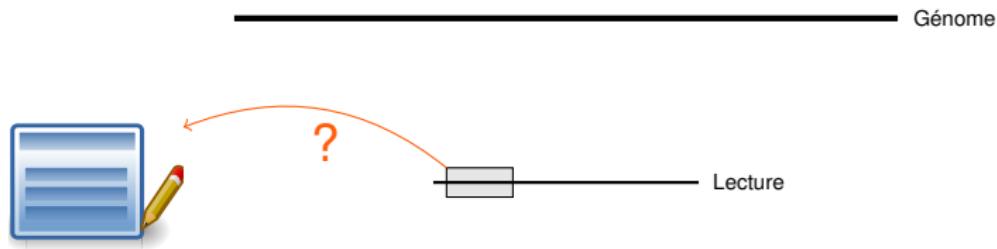
CC BY SA RRZEicons Wikimedia

Graines



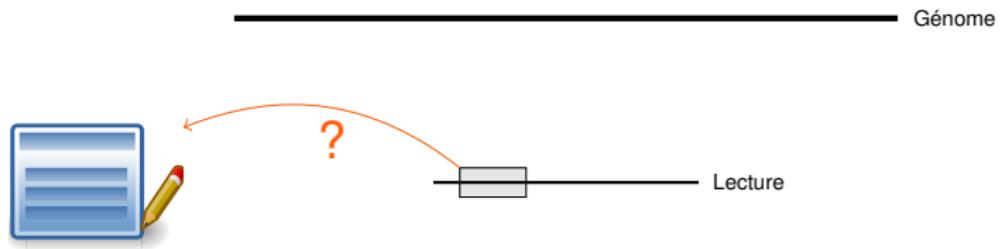
CC BY SA RRZEicons Wikimedia

Graines



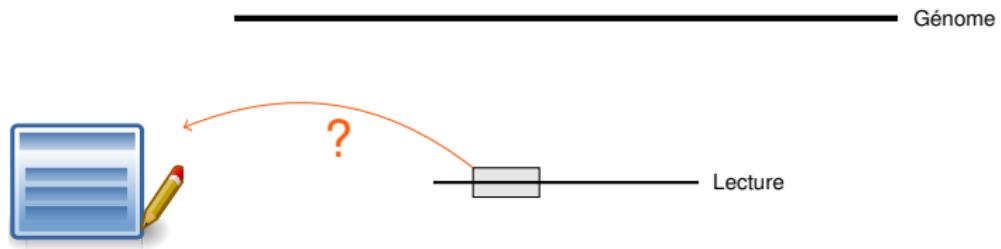
CC BY SA RRZEicons Wikimedia

Graines



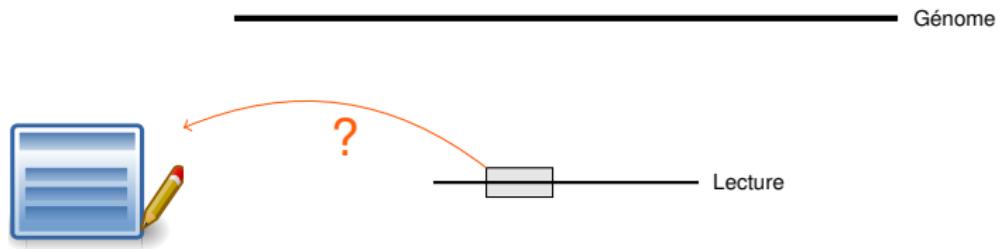
CC BY SA RRZEicons Wikimedia

Graines



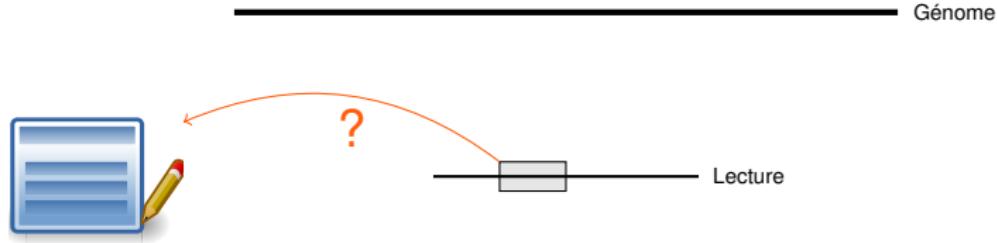
CC BY SA RRZEicons Wikimedia

Graines



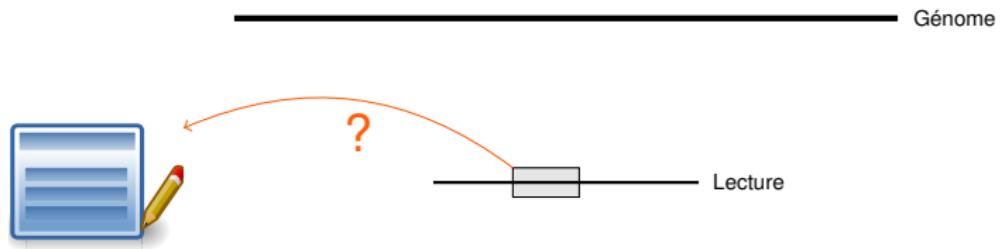
CC BY SA RRZEicons Wikimedia

Graines



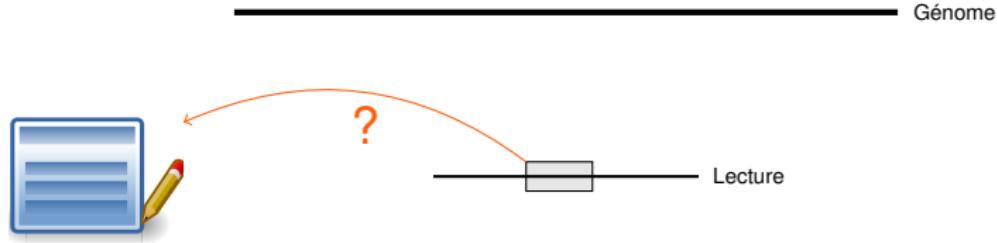
CC BY SA RRZEicons Wikimedia

Graines



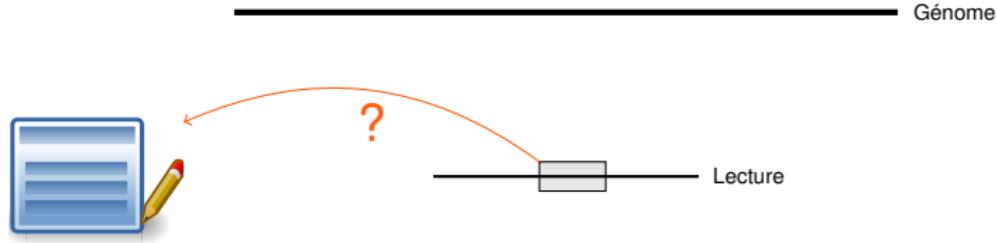
CC BY SA RRZEicons Wikimedia

Graines



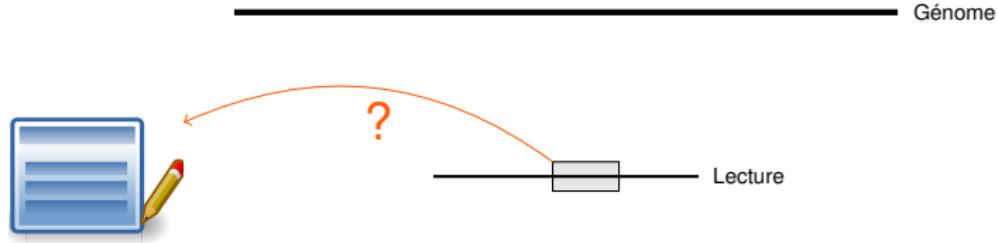
CC BY SA RRZEicons Wikimedia

Graines



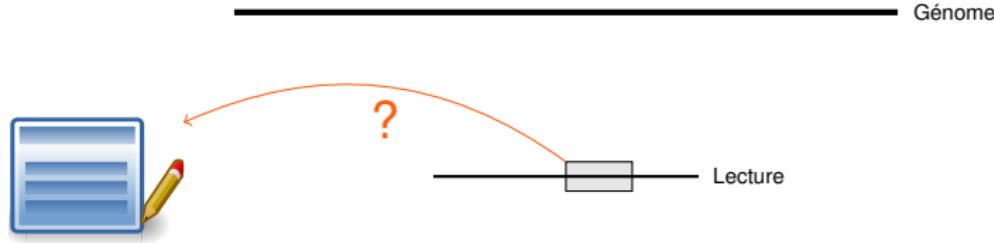
CC BY SA RRZEIcons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines

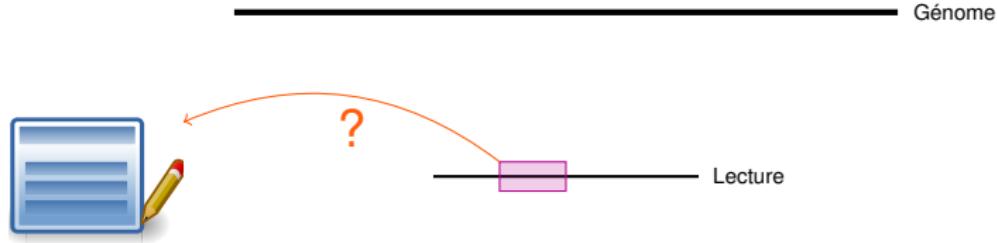
— Génome



— Lecture

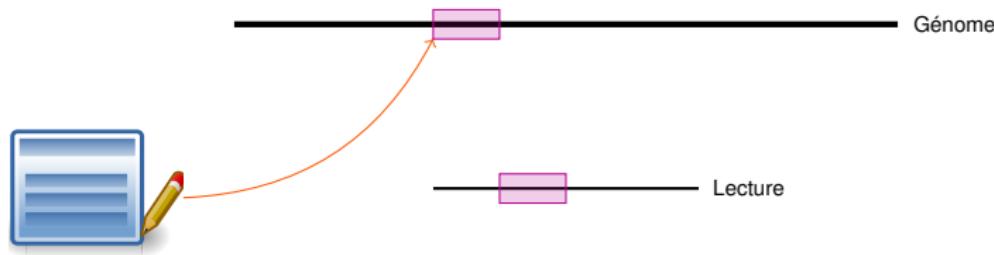
CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEicons Wikimedia

Graines



Génome



Lecture



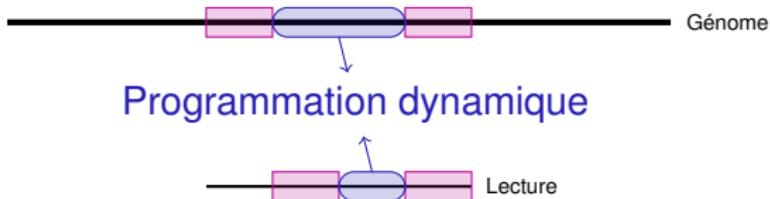
CC BY SA RRZEicons Wikimedia

Graines



CC BY SA RRZEIcons Wikimedia

Graines



CC BY SA RRZEIcons Wikimedia

Graines



Génome



Lecture

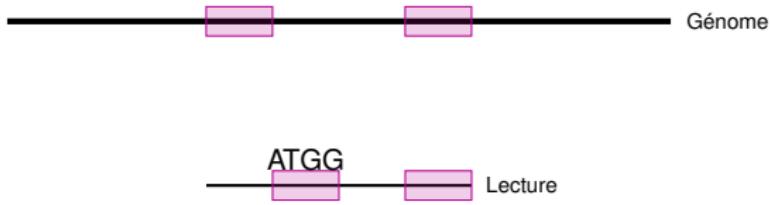


CC BY SA RRZEicons Wikimedia

Graines contiguës

####

Graines

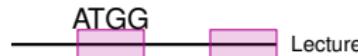


CC BY SA RRZEicons Wikimedia

Graines contiguës

####

Graines



CC BY SA RRZEicons Wikimedia

Graines contiguës

####

Graines

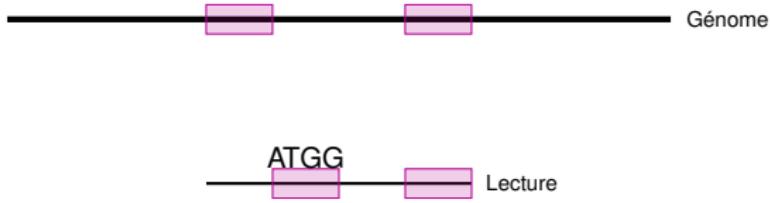


CC BY SA RRZEicons Wikimedia

Graines espacées

##-#

Graines

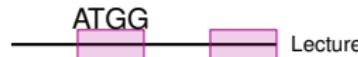


CC BY SA RRZEicons Wikimedia

Graines espacées

##-#

Graines

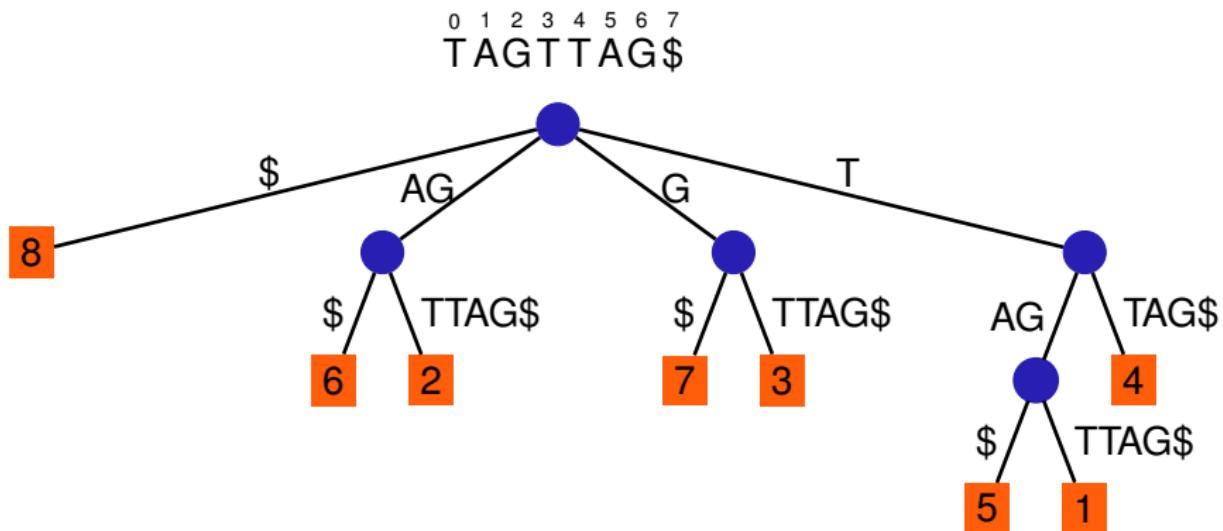


CC BY SA RRZEIcons Wikimedia

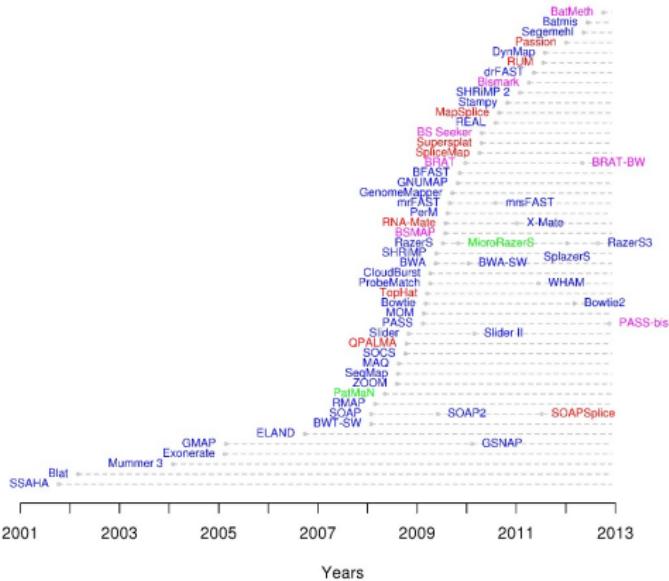
Graines espacées

##-#

Indexation de tout le texte



Les aligneurs



© Fonseca et al Bioinformatics, 2012

http://wwwdev.ebi.ac.uk/fg/hts_mappers/

Les aligneurs les plus utilisés

Bowtie

Blat

MAQ

BWA

TopHat

SOAP

SOAP2

Mummer3

BWA-SW

mrFAST

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al.*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés

Bowtie

Blat

MAQ

BWA

TopHat

SOAP

SOAP2

Mummer3

BWA-SW

mrFAST

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al.*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés

Bowtie

Blat

BWT

MAQ

BWA

TopHat

SOAP

SOAP2

Mummer3

BWA-SW

mrFAST

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés

Bowtie

Blat

BWT

MAQ

BWA

TopHat

Arbre des suffixes

SOAP

SOAP2

Mummer3

BWA-SW

mrFAST

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés

Bowtie

Blat

BTW

MAQ

BWA

TopHat

Arbre des suffixes

SOAP

SOAP2

Graines contiguës

Mummer3

BWA-SW

mrFAST

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés

Bowtie

Blat

BTW

MAQ

BWA

TopHat

Arbre des suffixes

SOAP

SOAP2

Graines contiguës

Mummer3

BWA-SW

mrFAST

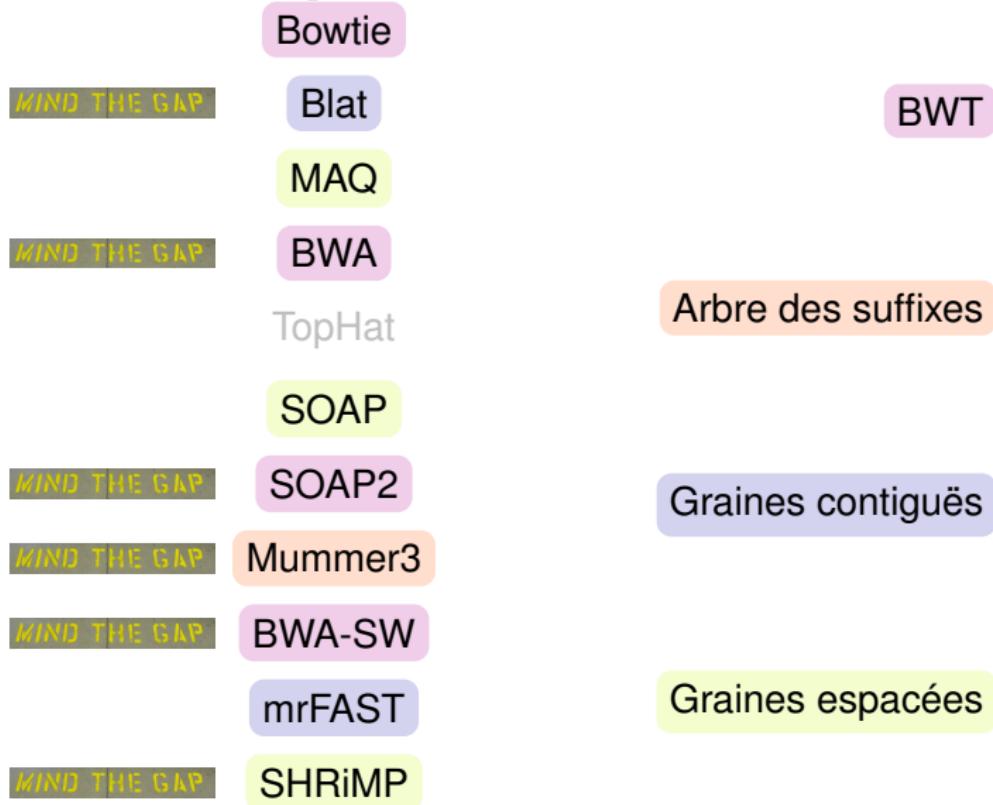
Graines espacées

SHRiMP

Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al.*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Les aligneurs les plus utilisés



Les aligneurs ayant le plus de citations par an, d'après Fonseca *et al*, Bioinformatics 2012.

Alignement, la spécificité d'abord ! — mikael.salson@lifl.fr

Qu'est-ce qu'un bon alignement ?

Quelle utilisation ?

Qu'est-ce qu'un bon alignement ?

Quelle utilisation ?

Quel protocole ?

Qu'est-ce qu'un bon alignement ?

Quelle utilisation ?

Quel protocole ?

Quel séquenceur ?

Qu'est-ce qu'un bon alignement ?

Quelle utilisation ?

Quel protocole ?

Quel séquenceur ?

Quelles machines ?

Qu'est-ce qu'un bon alignement ?

Quelle utilisation ?

Quel protocole ?

Tout est relatif !

Quel séquenceur ?

Quelles machines ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Génomé ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Génomé ?

Quantité ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Génomé ?

Quantité ?

Paires ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Génomé ?

Quantité ?

Paires ?

Erreur ?

Qu'est-ce qu'un bon alignement ?

Pourcentage de lectures correctement alignées ?

Longueur ?

Génomé ?

Quantité ?

Paires ?

Erreur ?

Alignement correct ?

La qualité, une donnée fiable ?

$$Q = -10 \log_{10} p$$

La qualité, une donnée fiable ?

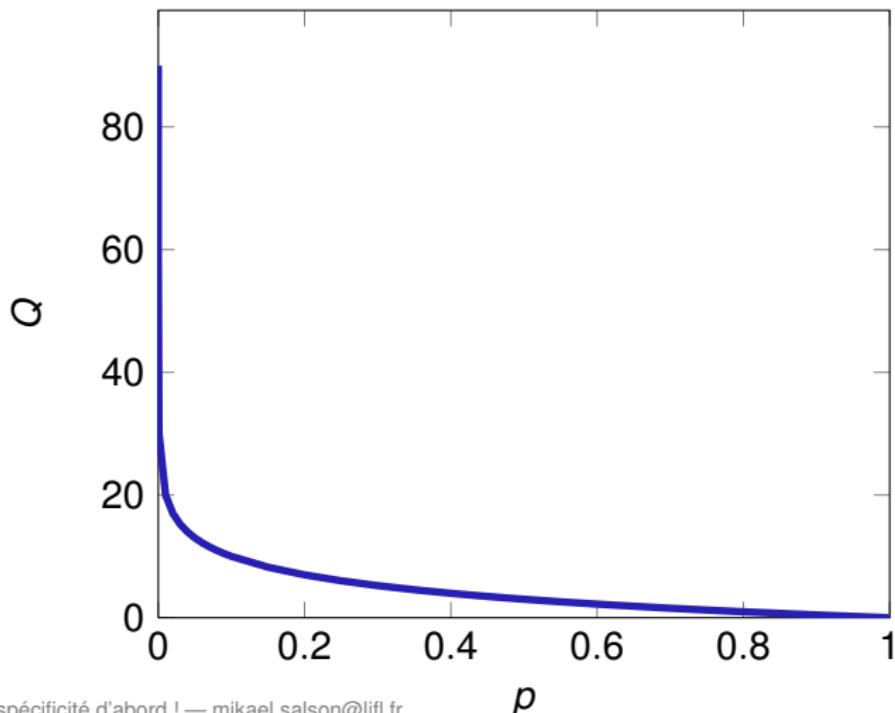
$$Q = -10 \log_{10} p$$

probabilité d'avoir
une erreur

La qualité, une donnée fiable ?

$$Q = -10 \log_{10} p$$

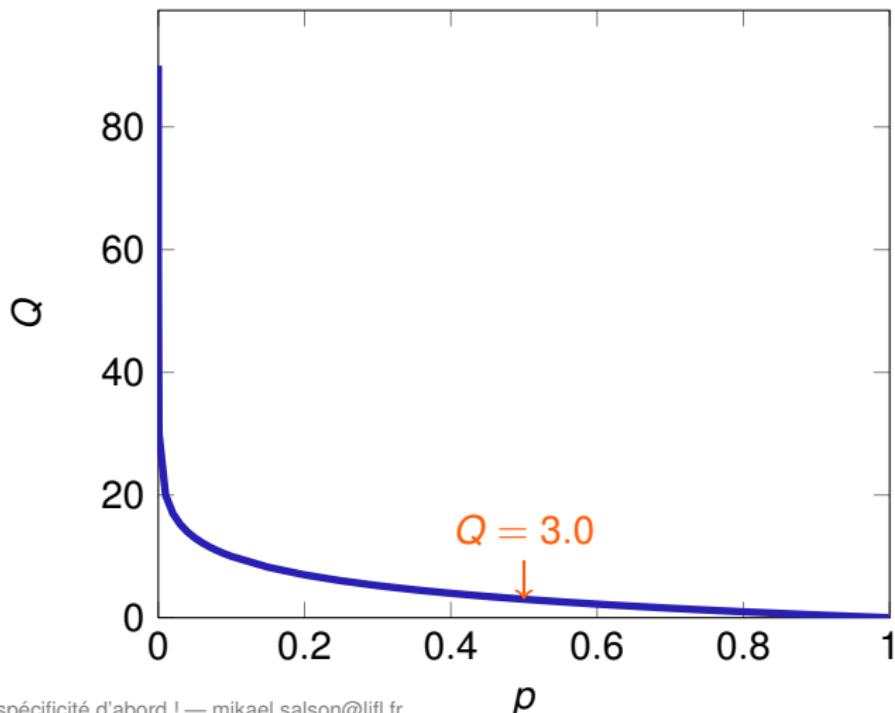
probabilité d'avoir
une erreur



La qualité, une donnée fiable ?

$$Q = -10 \log_{10} p$$

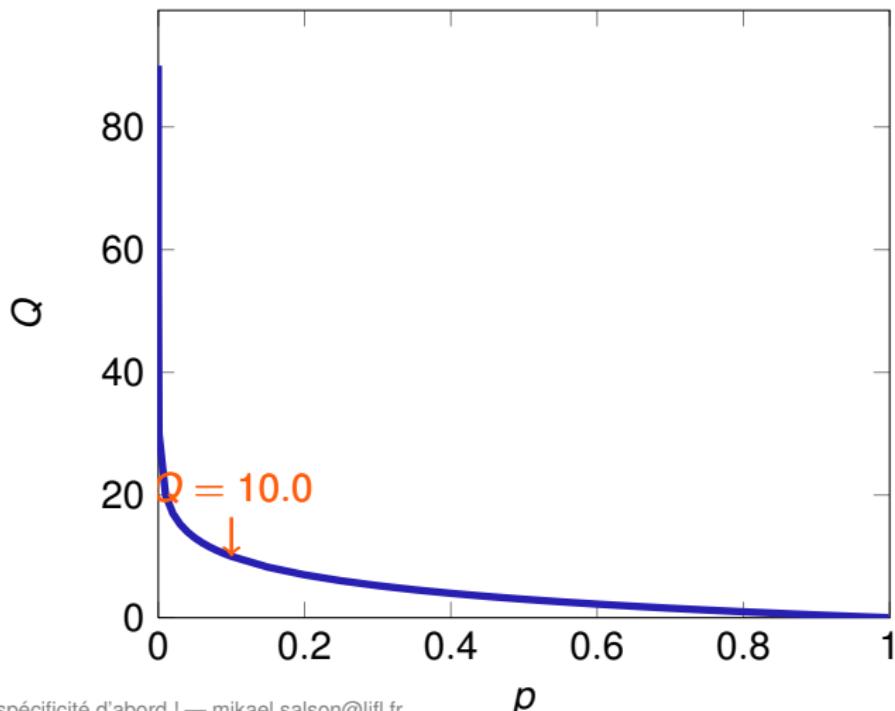
probabilité d'avoir
une erreur



La qualité, une donnée fiable ?

$$Q = -10 \log_{10} p$$

probabilité d'avoir
une erreur



La qualité, une donnée fiable !

La qualité, une donnée fiable !

« Our results indicate that significant gains in Solexa read mapping performance can be achieved by [...] appropriately using the base-call quality scores »

Smith *et al*, BMC Bioinformatics, 2008

La qualité, une donnée fiable !

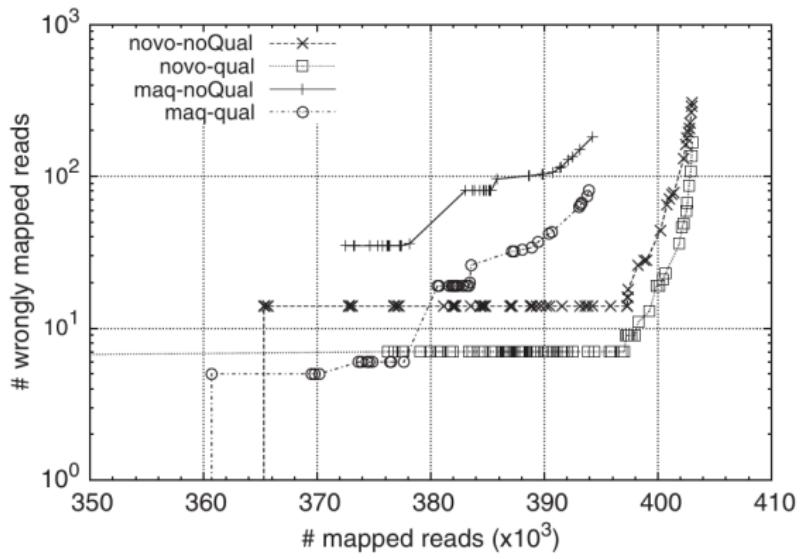
« Our results indicate that significant gains in Solexa read mapping performance can be achieved by [...] appropriately using the base-call quality scores »

Smith et al, BMC Bioinformatics, 2008

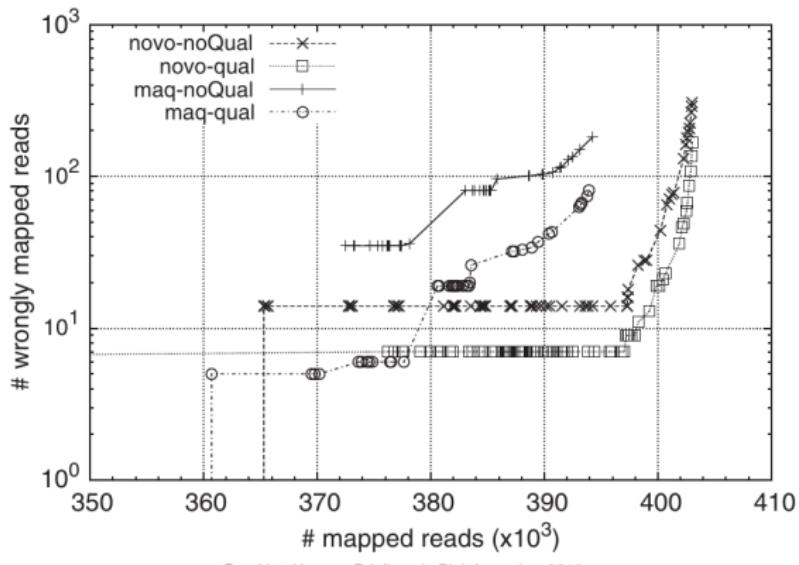
Alignement avec
X substitutions

Alignement sans
substitution ou
permise si $Q \leq 8$

La qualité, une donnée fiable !



La qualité, une donnée fiable !



© Li et Homer Briefings in Bioinformatics, 2010

« *Figure 3 shows that using base quality score halves alignment errors when the quality score is accurate* »

Li et Homer, Briefings in Bioinformatics, 2010

La qualité, une donnée fiable ! ?

La qualité, une donnée fiable ! ?

chr. 21
humain

DRX000307

Illumina GA II
76 pb
20 M lectures

La qualité, une donnée fiable ! ?

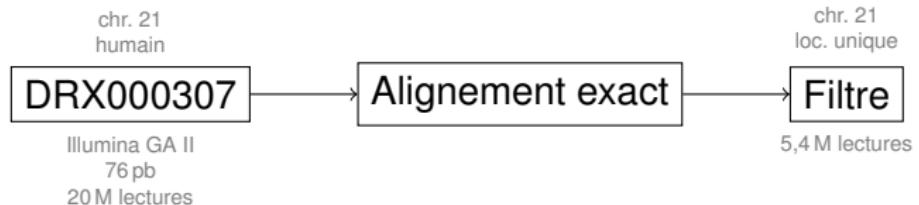
chr. 21
humain

DRX000307

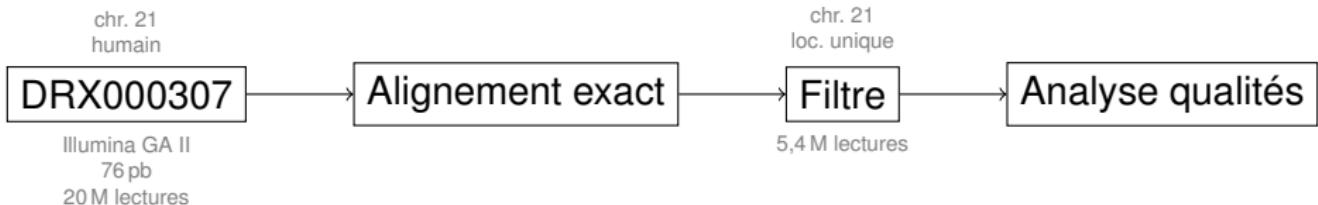
Alignement exact

Illumina GA II
76 pb
20 M lectures

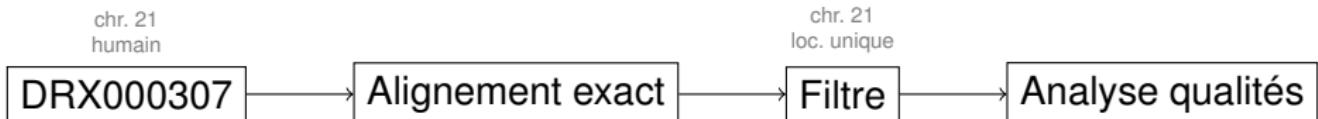
La qualité, une donnée fiable ! ?



La qualité, une donnée fiable ! ?

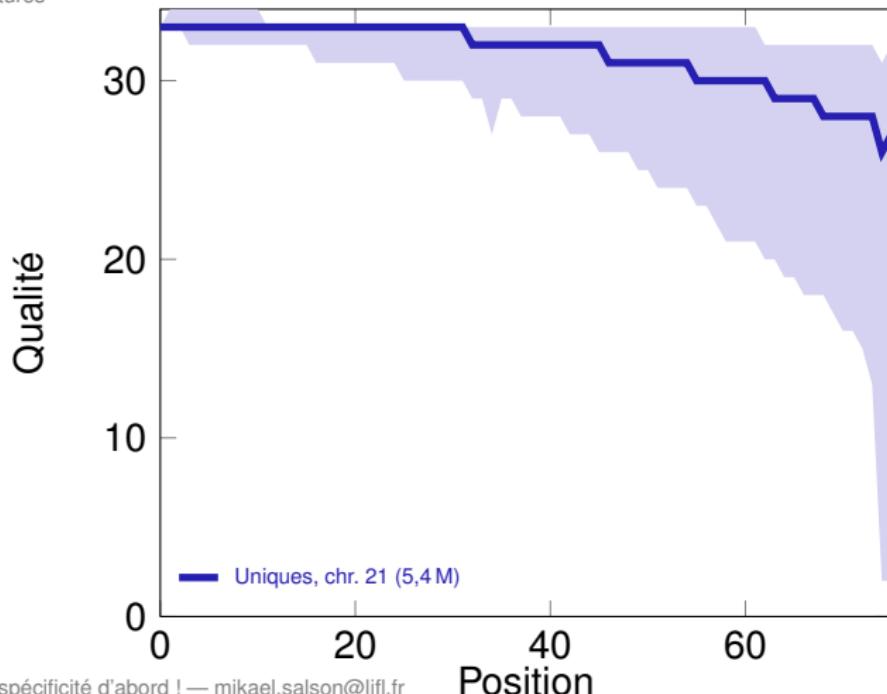


La qualité, une donnée fiable ! ?

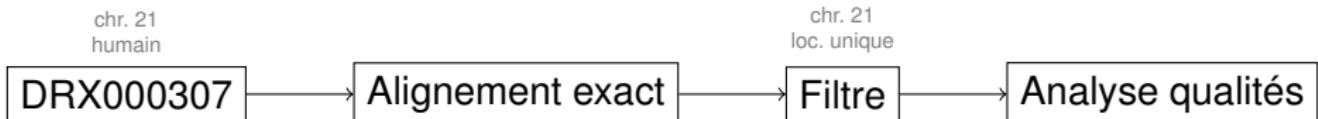


Illumina GA II
76 pb
20 M lectures

5,4 M lectures

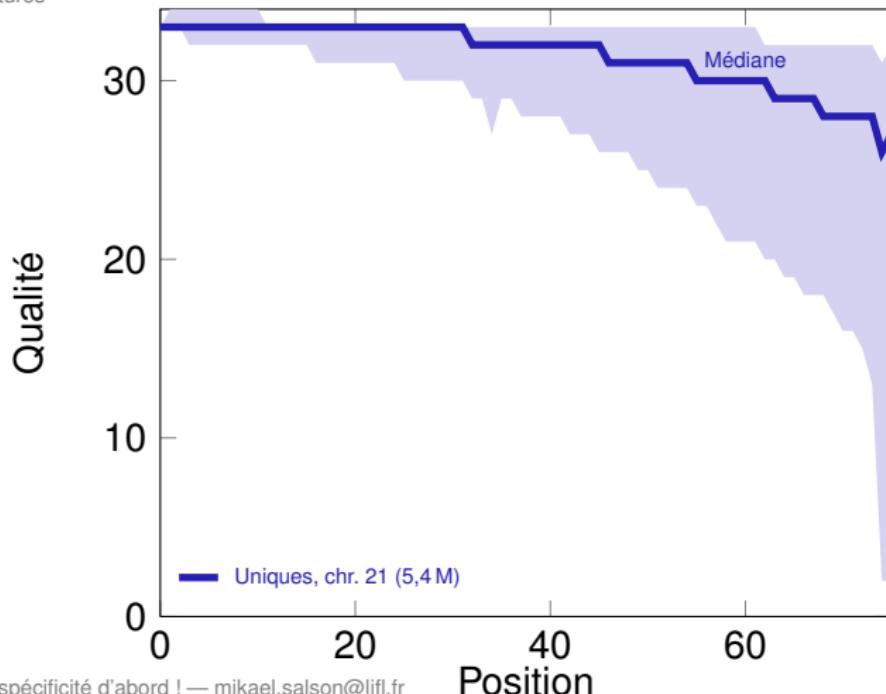


La qualité, une donnée fiable ! ?



Illumina GA II
76 pb
20 M lectures

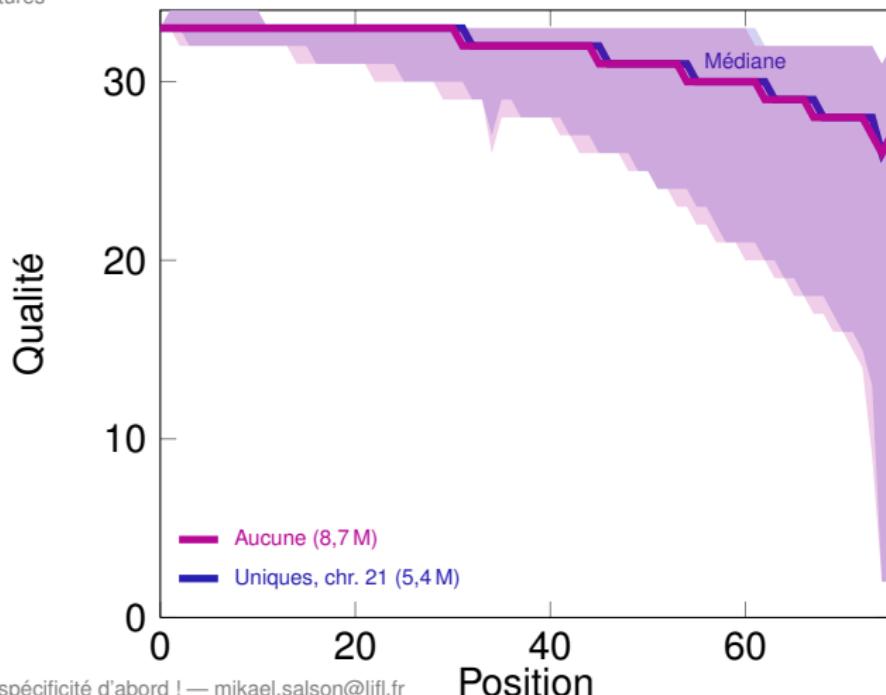
5,4 M lectures



La qualité, une donnée fiable ! ?



Illumina GA II
76 pb
20 M lectures



Qualité, une donnée fiable ! ? ?

Qualité, une donnée fiable ! ? ?

chr. 21
humain

DRX000307

Illumina GA II
76 pb
20 M lectures

Qualité, une donnée fiable ! ? ?

chr. 21
humain

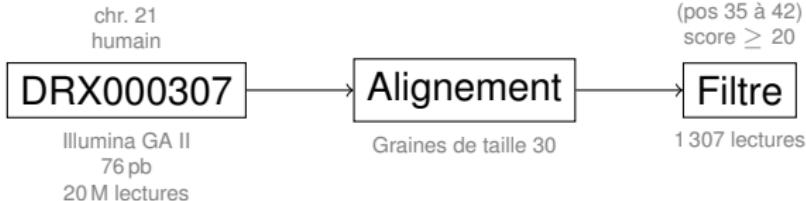
DRX000307

Alignement

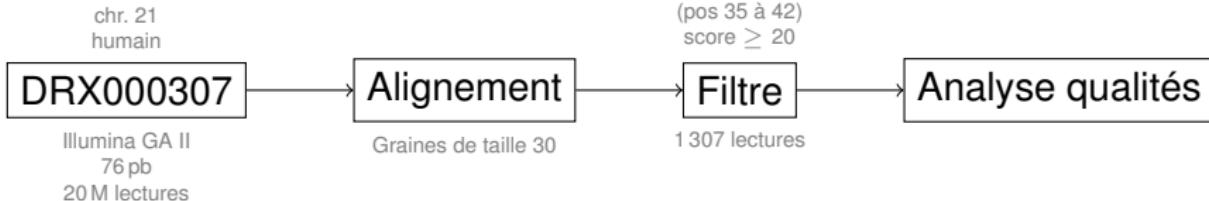
Illumina GA II
76 pb
20 M lectures

Graines de taille 30

Qualité, une donnée fiable ! ? ?



Qualité, une donnée fiable ! ? ?

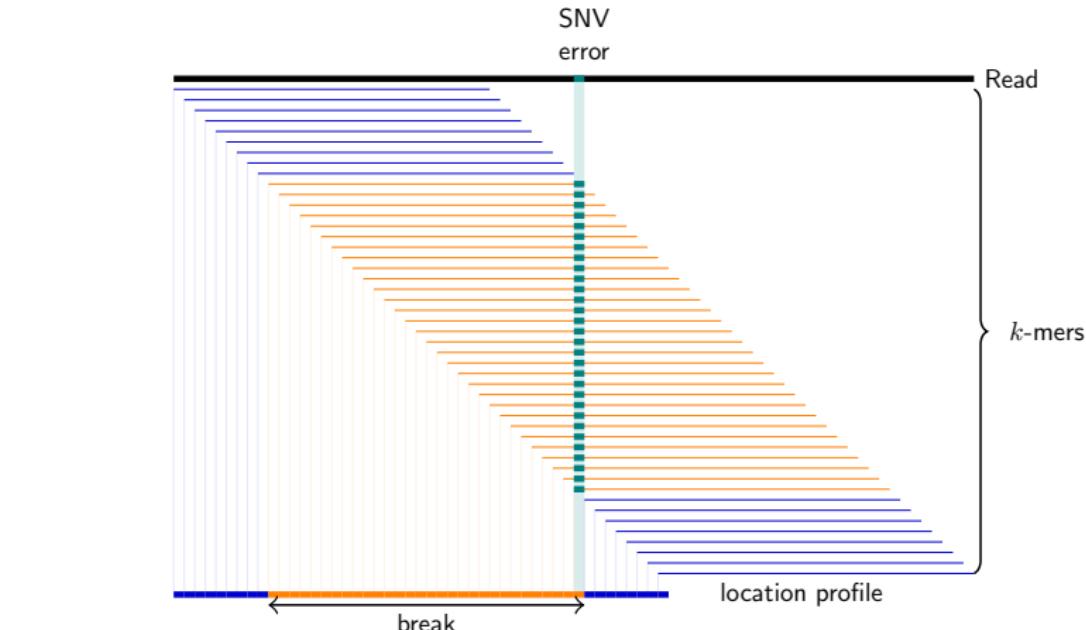


Qualité, une donnée fiable ! ? ?

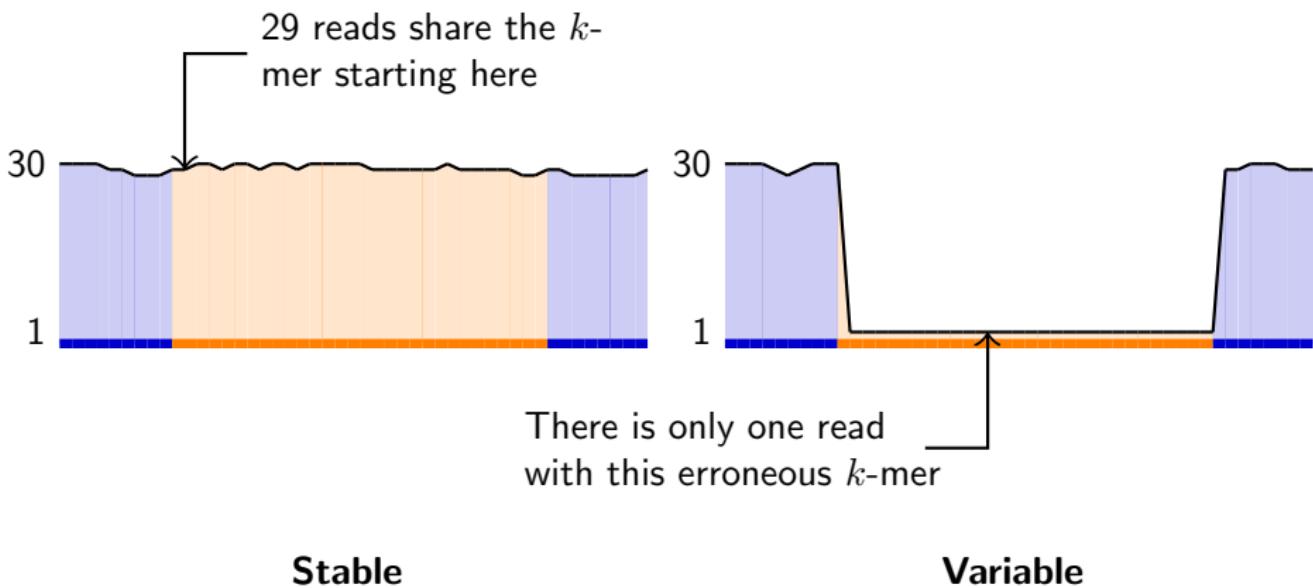
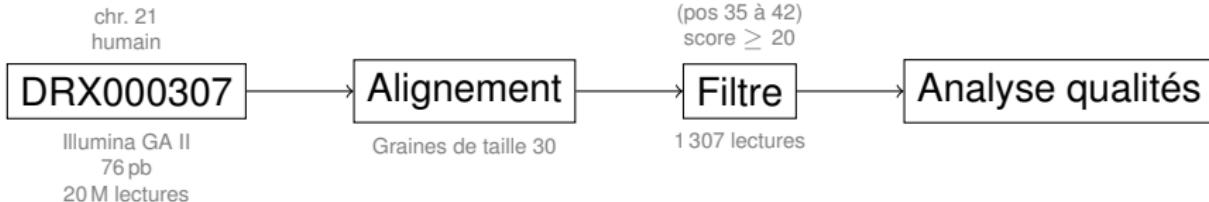


Illumina GA II
76 pb
20 M lectures

chr. 21
Substitution
(pos 35 à 42)
score ≥ 20



Qualité, une donnée fiable ! ? ?



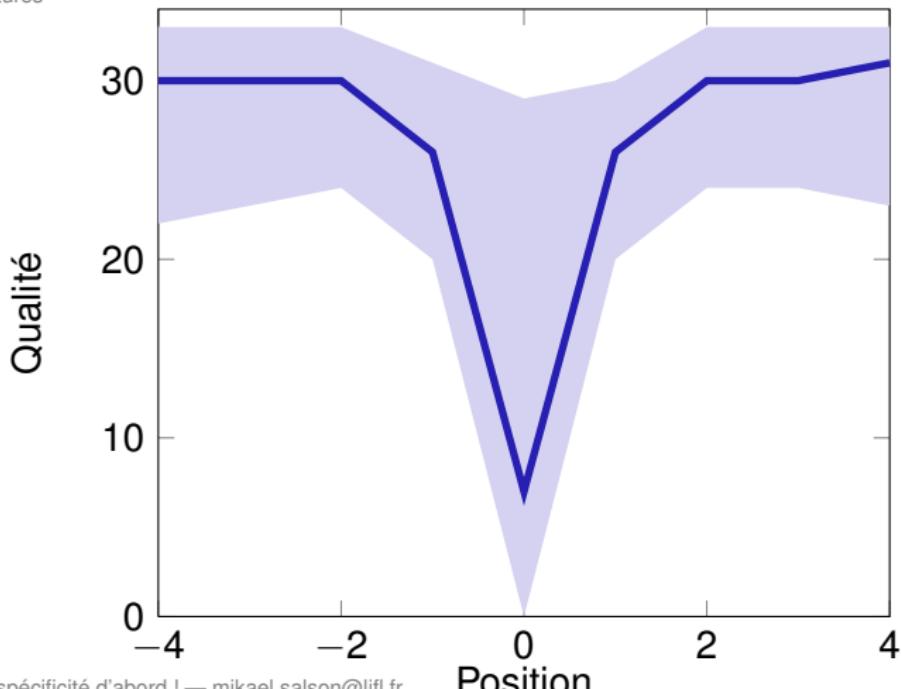
Qualité, une donnée fiable ! ? ?



Illumina GA II
76 pb
20 M lectures

Graines de taille 30

1 307 lectures



Évaluation des performances

Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011

Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011

Génome

de référence
ou artificiel

Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011



50 pb, 500 000 lectures
Substitutions
Indels

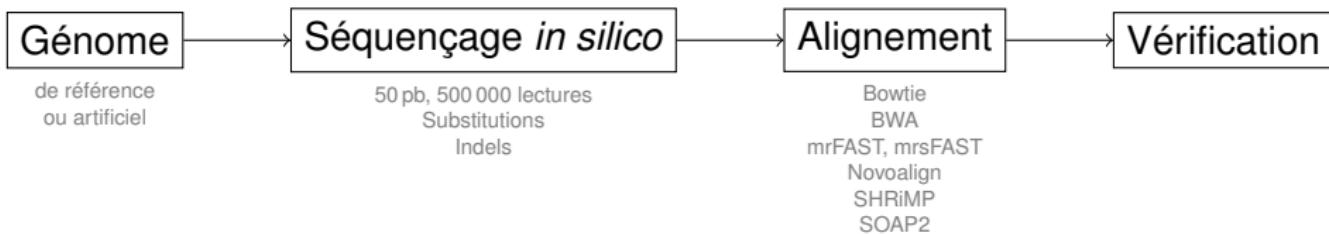
Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011



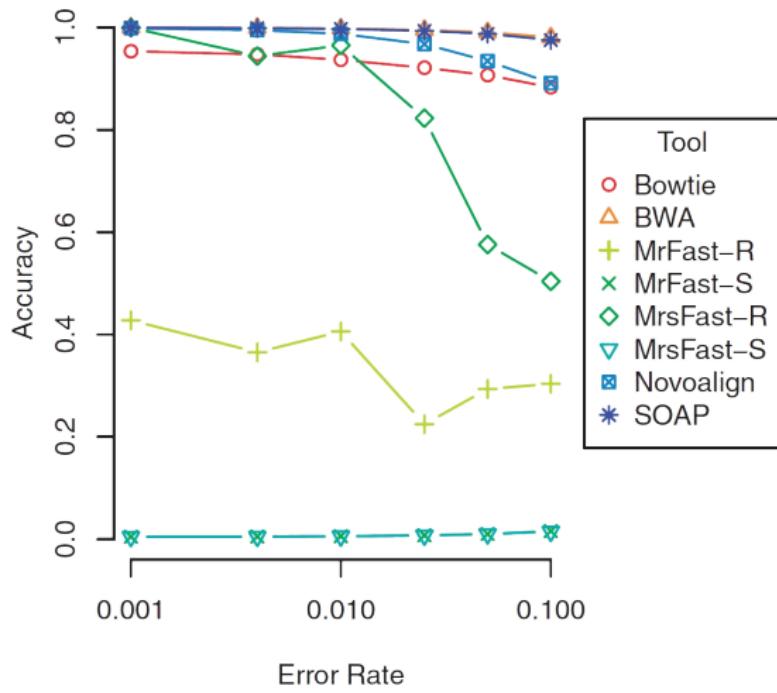
Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011



Évaluation des performances SEAL

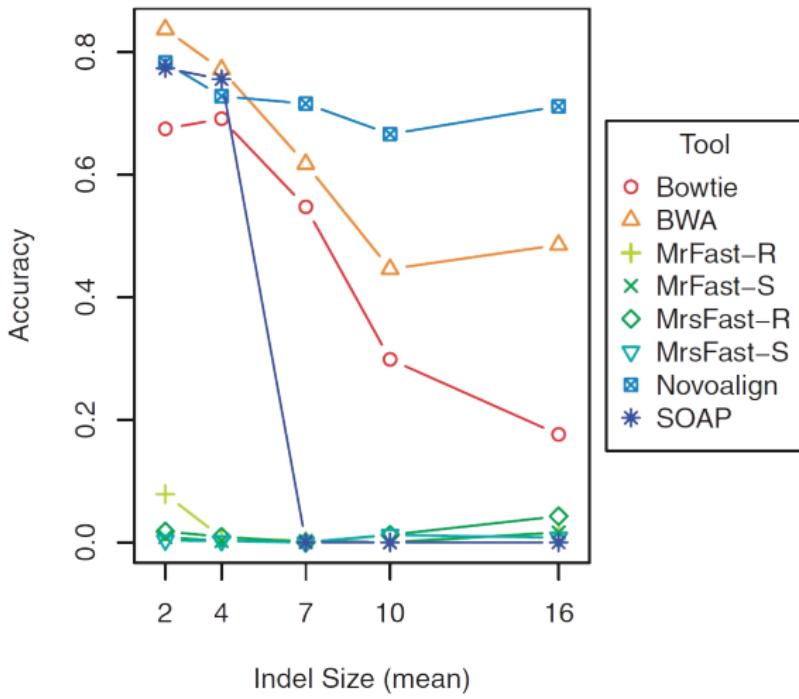
Ruffalo *et al*
Bioinformatics, 2011



© Ruffalo *et al* Bioinformatics

Évaluation des performances SEAL

Ruffalo *et al*
Bioinformatics, 2011



© Ruffalo *et al* Bioinformatics

Évaluation des performances Rabema

Holtgrewe *et al*
BMC Bioinformatics, 2011

Évaluation des performances

Rabema

Holtgrewe *et al*
BMC Bioinformatics, 2011

Données réelles

Évaluation des performances

Rabema

Holtgrewe *et al*
BMC Bioinformatics, 2011



10 000 lectures sélectionnées
Sensibilité max (RazerS)

Évaluation des performances

Rabema

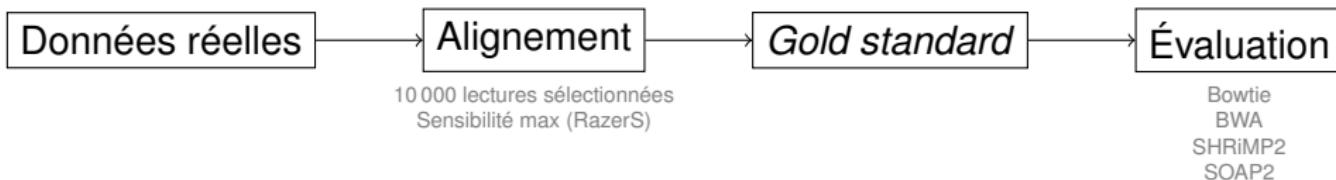
Holtgrewe *et al*
BMC Bioinformatics, 2011



Évaluation des performances

Rabema

Holtgrewe *et al*
BMC Bioinformatics, 2011

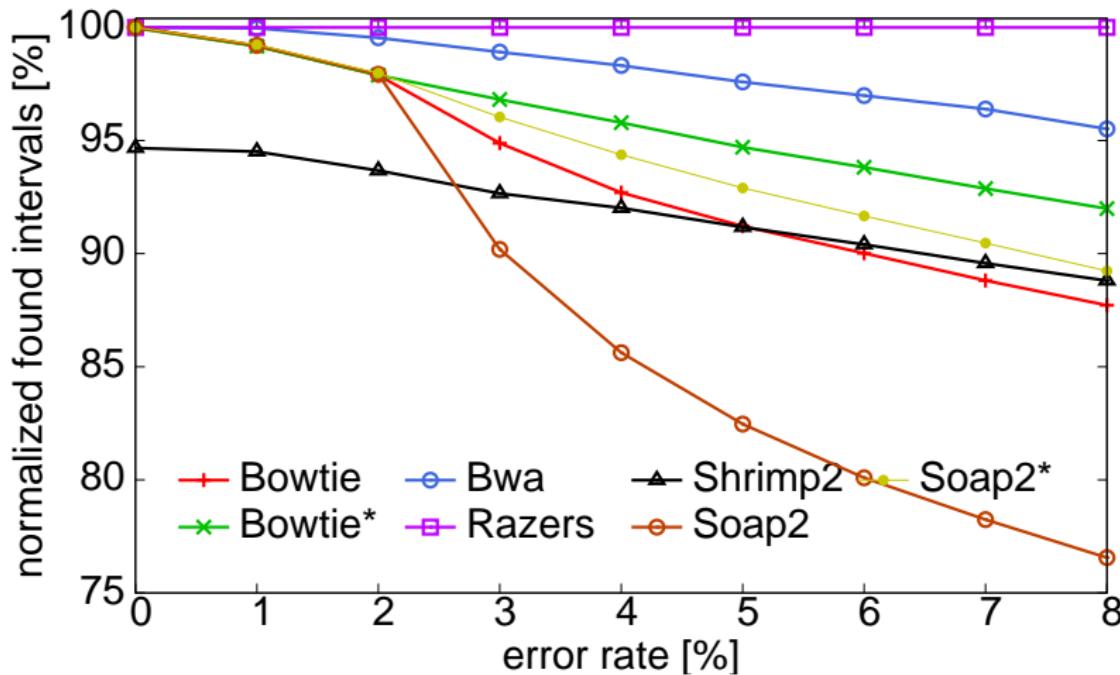


Évaluation des performances

Rabema

Holtgrewe *et al*
BMC Bioinformatics, 2011

D. melanogaster, lectures 100 pb Illumina



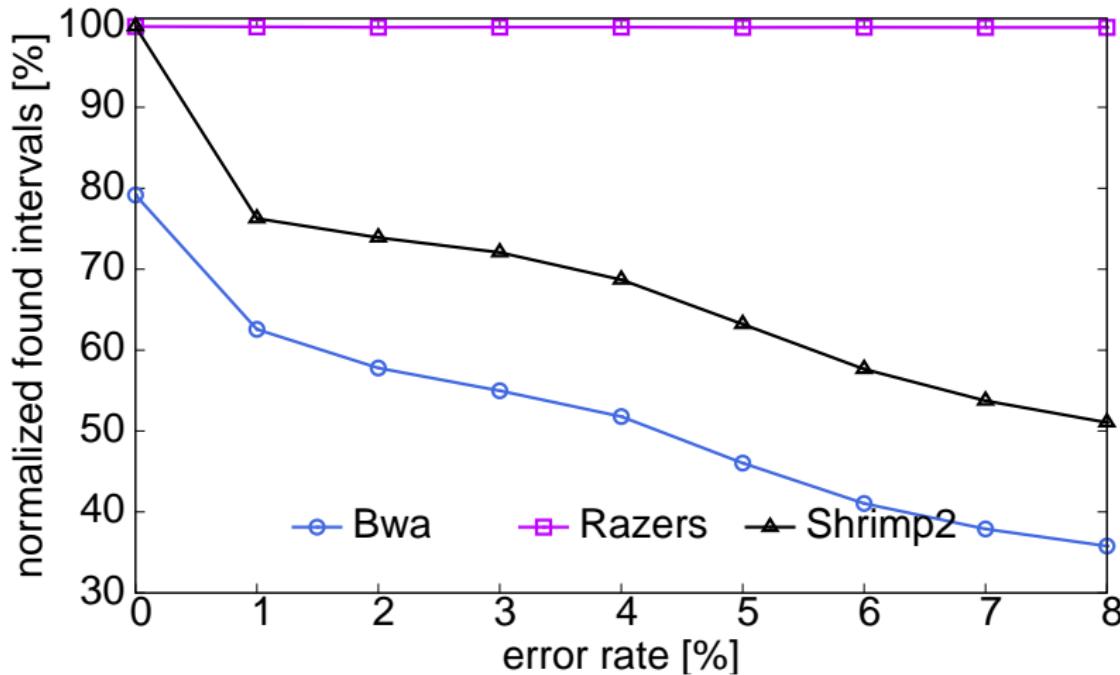
CC BY Holtgrewe *et al* BMC Bioinformatics

Évaluation des performances

Rabema

Holtgrewe et al
BMC Bioinformatics, 2011

D. melanogaster, lectures 454 (moyenne : 273 pb)



CC BY Holtgrewe et al BMC Bioinformatics

Évaluation des performances

Schbath *et al*

Journal of Computational Biology, 2012

Évaluation des performances

Schbath *et al*

Journal of Computational Biology, 2012

Génomes de
référence

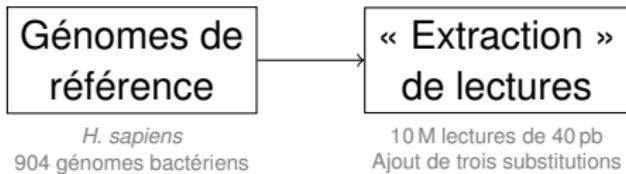
H. sapiens

904 génomes bactériens

Évaluation des performances

Schbath et al

Journal of Computational Biology, 2012



Évaluation des performances

Schbath et al

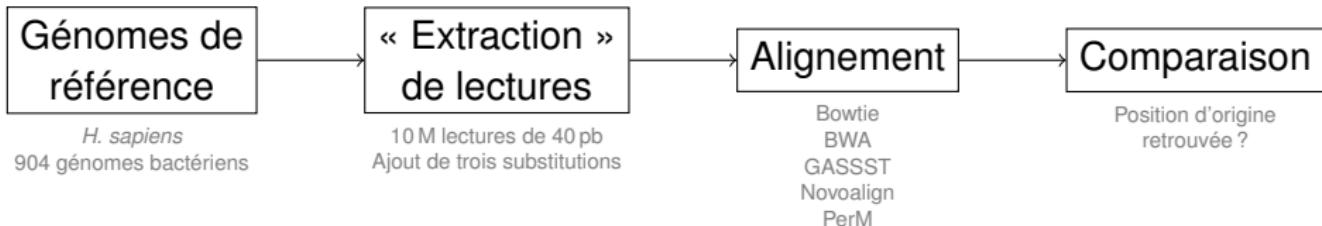
Journal of Computational Biology, 2012



Évaluation des performances

Schbath *et al*

Journal of Computational Biology, 2012

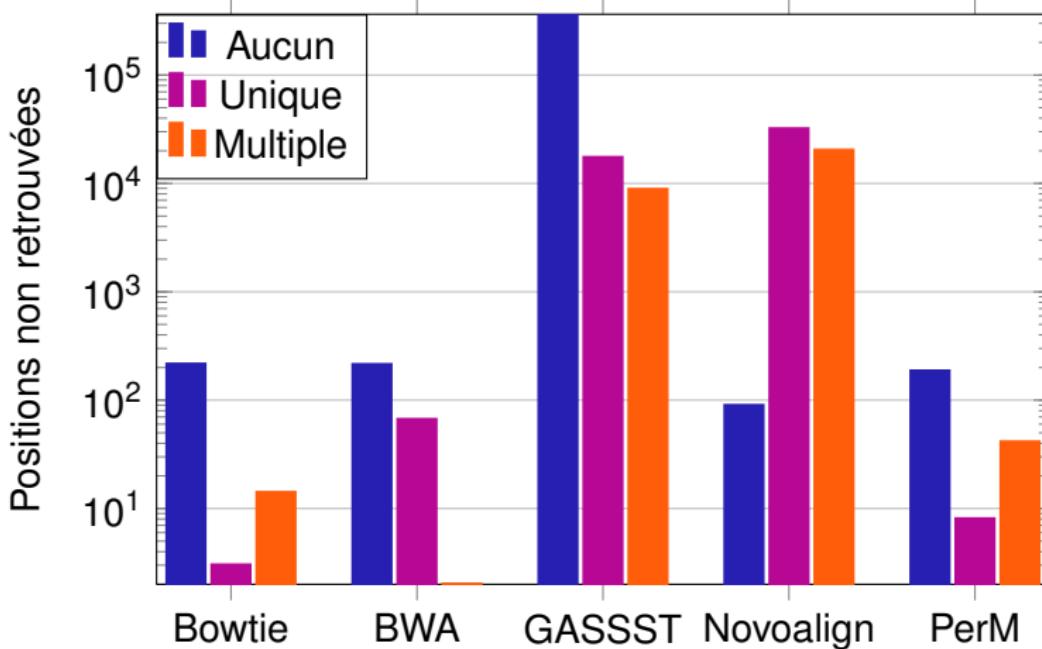


Évaluation des performances

Schbath *et al*

Journal of Computational Biology, 2012

Bactérien

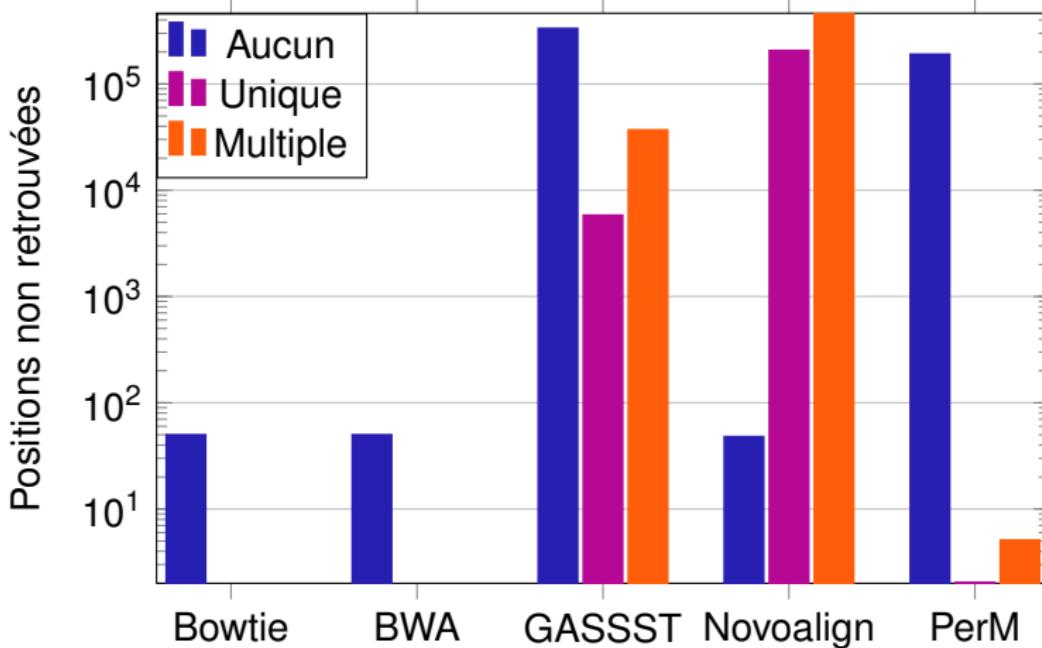


Évaluation des performances

Schbath *et al*

Journal of Computational Biology, 2012

Humain



Évaluation des performances

Yu et al

BioData Mining, 2012

Évaluation des performances

Yu et al

BioData Mining, 2012

Données réelles

3 000 exons humains
7,4 M et 5,4 M lectures de 68pb

Évaluation des performances

Yu et al

BioData Mining, 2012

Données réelles

3 000 exons humains
7,4 M et 5,4 M lectures de 68pb

→ Alignement

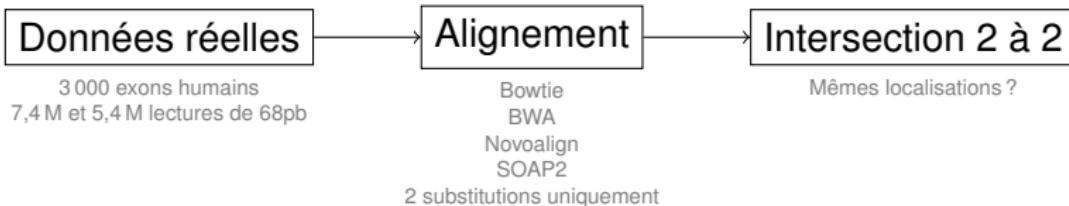
Bowtie
BWA
Novoalign
SOAP2

2 substitutions uniquement

Évaluation des performances

Yu et al

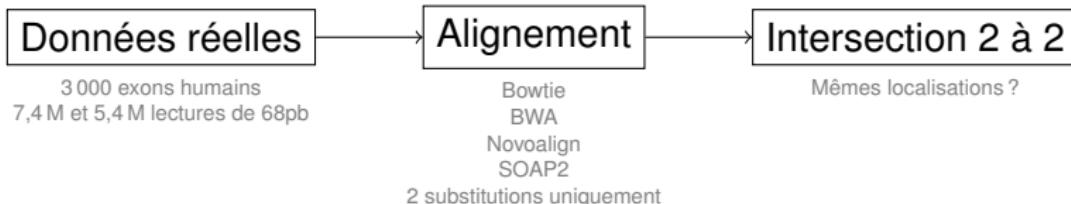
BioData Mining, 2012



Évaluation des performances

Yu et al

BioData Mining, 2012



~ 95 % d'identité entre Bowtie, BWA et SOAP2

~ 85 % pour Novoalign

Évaluation... et le RNA-Seq ?

Jonctions exon-exon ?

Évaluation... et le RNA-Seq ?

Jonctions exon-exon ?

Transcrits de fusion ?

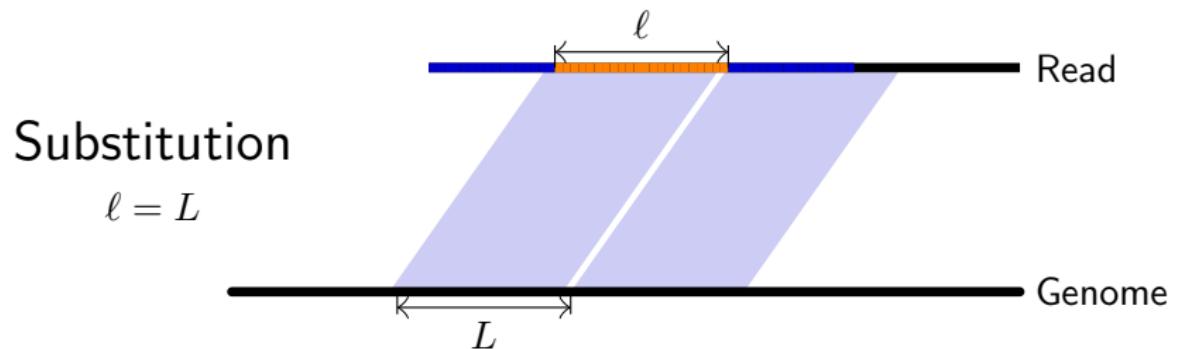
CRAC, un outil pour le RNA-Seq



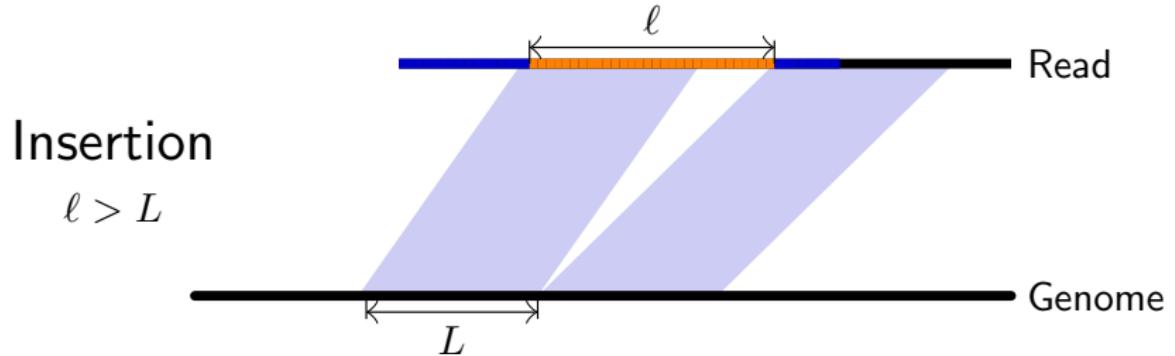
CRAC, un outil pour le RNA-Seq

Avec Nicolas Philippe, Thérèse Commes, Éric Rivals
crac.gforge.inria.fr

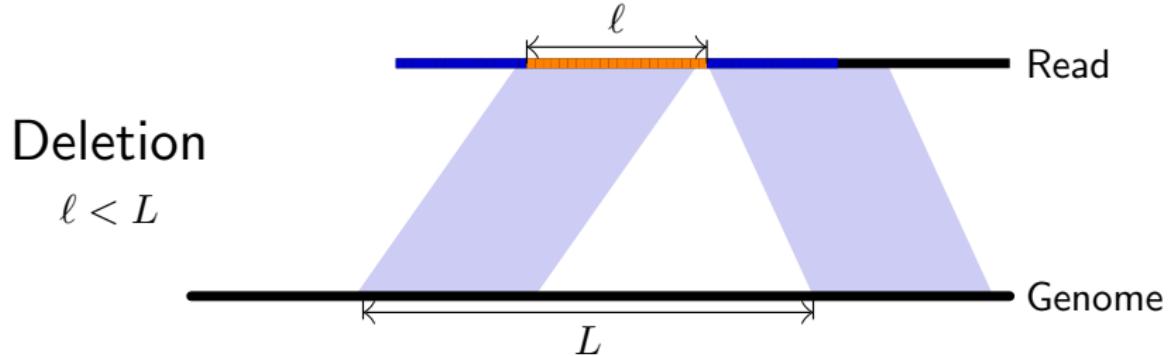
CRAC, un outil pour le RNA-Seq



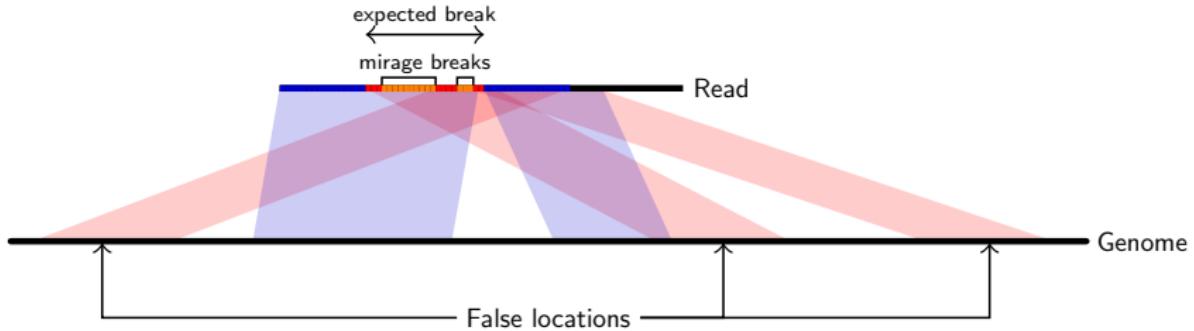
CRAC, un outil pour le RNA-Seq



CRAC, un outil pour le RNA-Seq



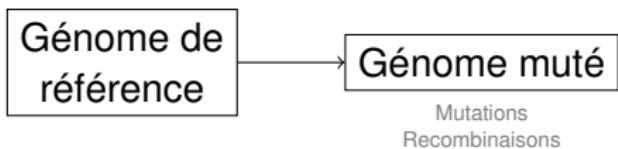
CRAC, un outil pour le RNA-Seq



Données simulées en RNA-Seq

Génome de
référence

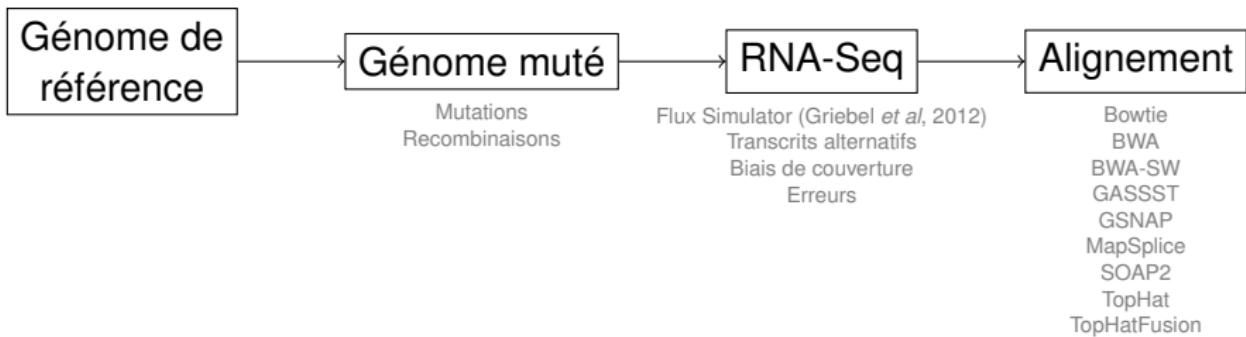
Données simulées en RNA-Seq



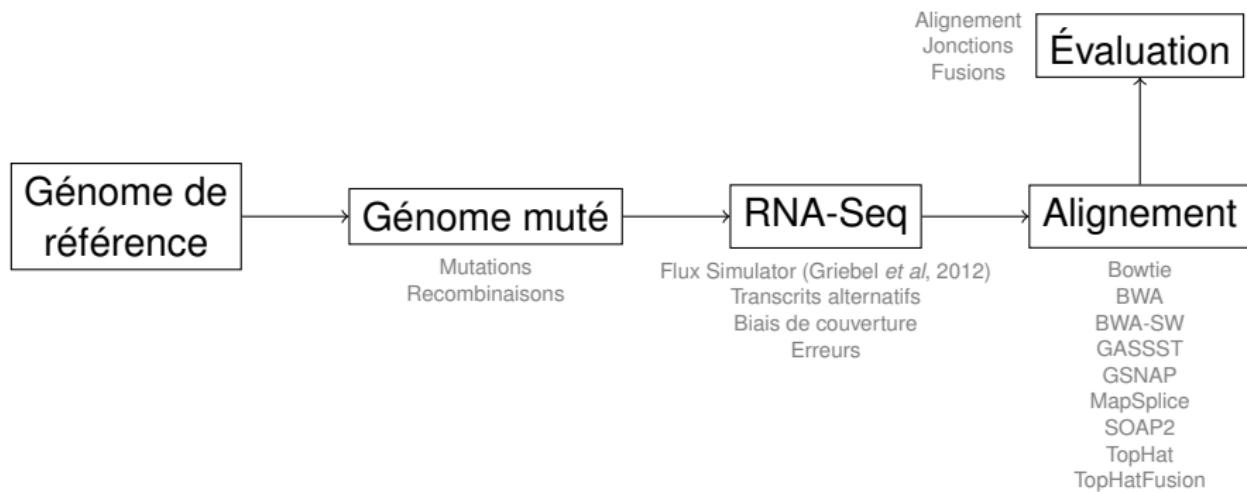
Données simulées en RNA-Seq



Données simulées en RNA-Seq



Données simulées en RNA-Seq



Données simulées

Type	SNV	Insertion	Délétion	Fusion	Erreurs
Génome complet	3 139 937	287 336	287 502	1 002	
45 M lectures (75 pb)	29 084	2 687	2 734	647	12 836 882
48 M lectures (200 pb)	52 971	4 810	4 901	914	38 840 045

Performances sur données simulées

Alignment (occurrences uniques)

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
Bowtie	75.42	99.59	55.72	99.81
BWA	79.29	99.13	68.66	96.86
CRAC	94.51	99.72	95.89	99.78
GASSST	70.73	99.09	59.43	97.86
GSNAP	94.62	99.88	84.84	99.28
SOAP2	77.6	99.52	56.08	99.78

Performances sur données simulées

Découverte de jonctions exon-exon

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
CRAC	79.43	99.5	86.02	99.18
GSNAP	84.17	97.03	72.94	97.09
MapSplice	79.89	97.68	84.72	98.82
TopHat	84.96	89.59	54.07	94.69

Performances sur données simulées

Découverte de transcrits de fusion

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
CRAC	53.89	93.84	64.86	90.18
MapSplice	2.33	0	2.63	0.01
TopHatFusion	32.73	42.02		
TopHatFusionPost	12.26	97.22		

Jonctions dans données réelles

ERR030856

RNA-Seq, Hi-Seq 2000

100 pb orienté

76 M lectures

Jonctions dans données réelles

ERR030856

RNA-Seq, Hi-Seq 2000
100 pb orienté
76 M lectures

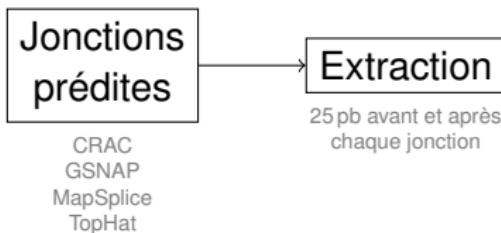
Jonctions
prédites

CRAC
GSNAP
MapSplice
TopHat

Jonctions dans données réelles

ERR030856

RNA-Seq, Hi-Seq 2000
100 pb orienté
76 M lectures



Jonctions dans données réelles

ERR030856

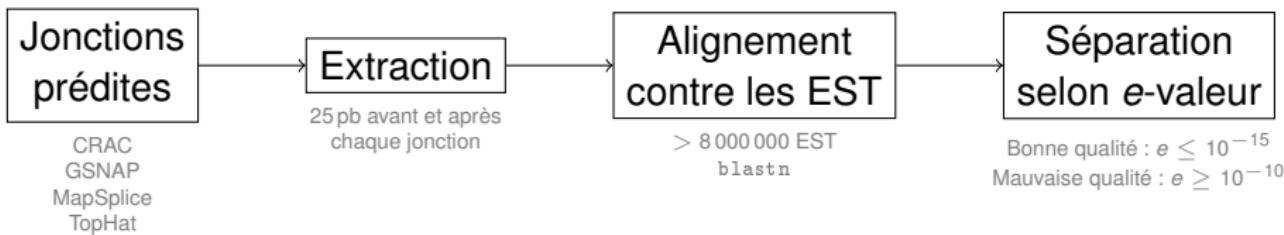
RNA-Seq, Hi-Seq 2000
100 pb orienté
76 M lectures



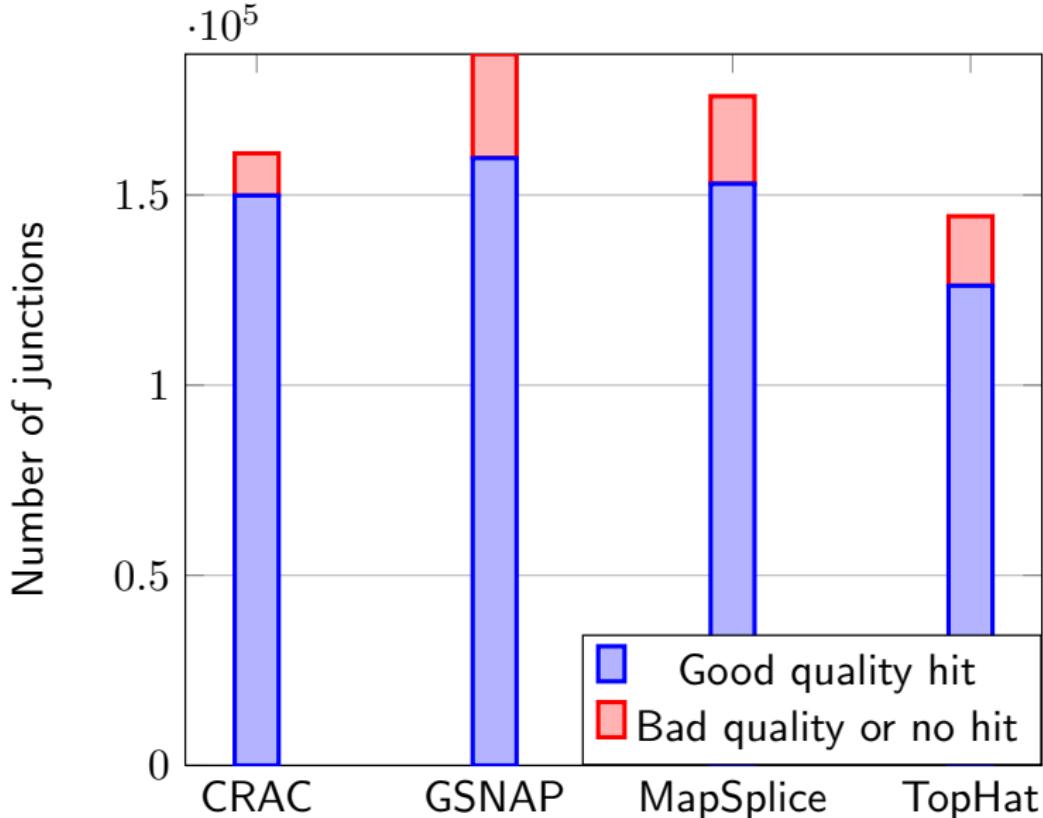
Jonctions dans données réelles

ERR030856

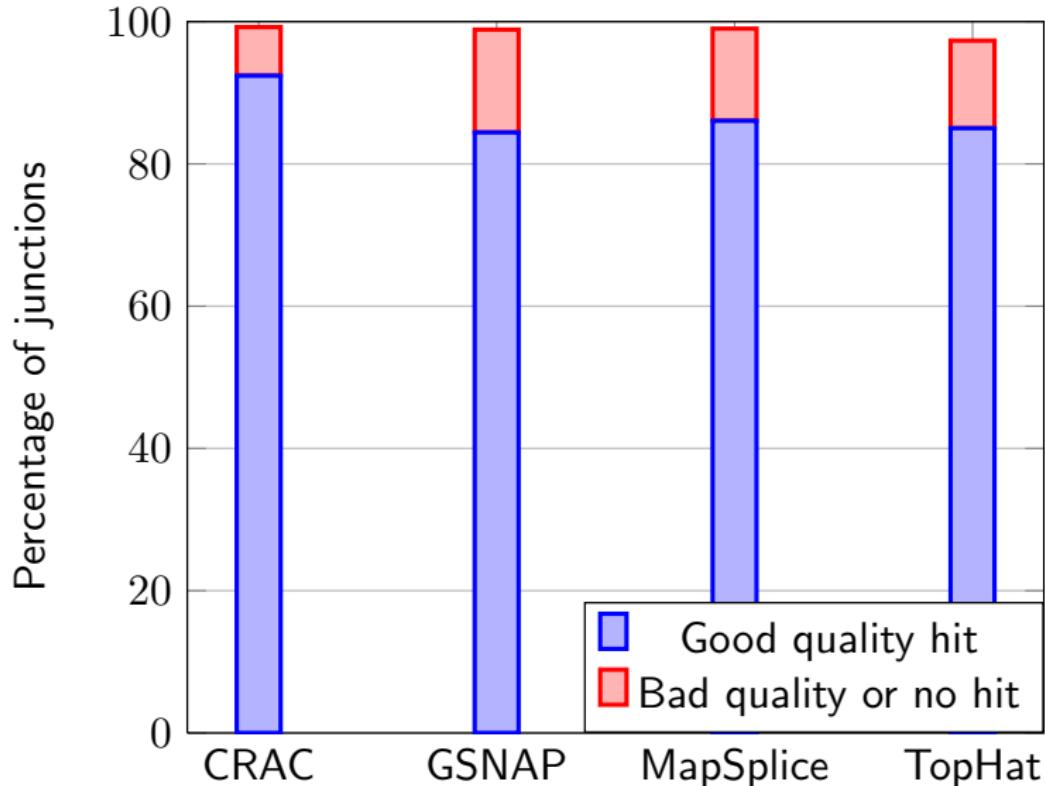
RNA-Seq, Hi-Seq 2000
100 pb orienté
76 M lectures



Jonctions dans données réelles



Jonctions dans données réelles



Fusions dans données réelles

Données réelles

Cellules de cancer du sein

27 fusions connues

Lectures 50 pb

Edgren *et al*, Genome Biology, 2011

Fusions dans données réelles



Cellules de cancer du sein

27 fusions connues

Lectures 50 pb

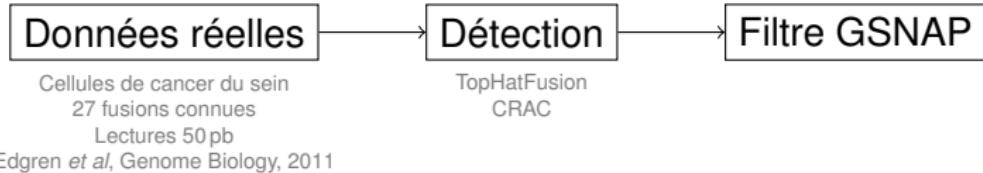
Edgren *et al*, Genome Biology, 2011

Détection

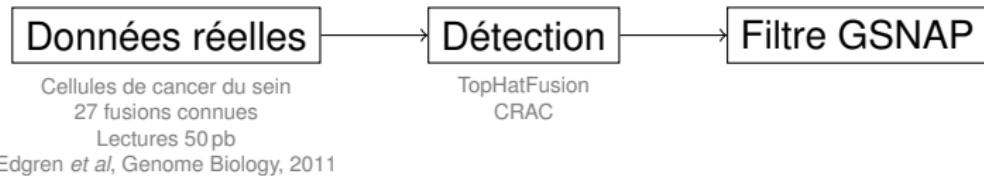
TopHatFusion

CRAC

Fusions dans données réelles

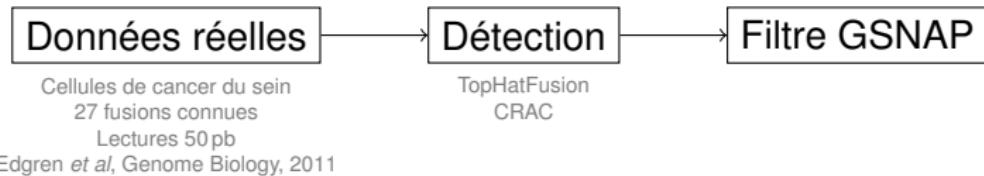


Fusions dans données réelles



CRAC
455

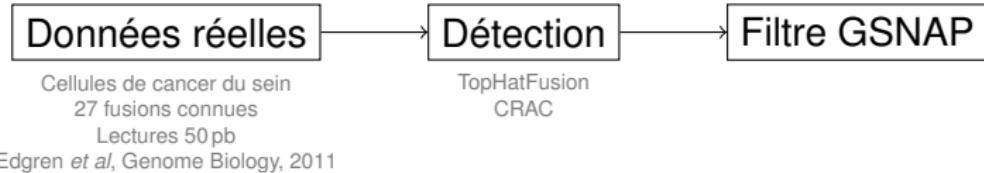
Fusions dans données réelles



CRAC
455

TopHatFusion
193 163

Fusions dans données réelles



CRAC

455

Validées : 20

TopHatFusion

193 163

Validées : 21

Perspectives...

Qualité : vraiment informative ?

Perspectives...

Qualité : vraiment informative ?

Nécessité d'évaluation plus complète

Bonus

Transformée de Burrows-Wheeler

$T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \text{CGAGACGAA\$} \end{matrix}$

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

C G A G A C G A A \\$

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

C G A G A C G A A \$
\$ C G A G A C G A A

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

C G A G A C G A A \$

\$ C G A G A C G A A

A \$ C G A G A C G A

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

C G A G A C G A A \$

\$ C G A G A C G A A

A \$ C G A G A C G A

A A \$ C G A G A C G

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA \$
\$ CGAGACGAA
A \$ CGAGACGA
AA \$ CGAGACG
GAA \$ CGAGAC
CGAA \$ CGAGA
ACGAA \$ CGAG
GACGAA \$ CGA
AGACGAA \$ CG
GAGACGAA \$ C

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$
\$CGAGACGAA
A\$CGAGACGA
AA\$CGAGACG
GAA\$CGAGAC
CGAA\$CGAGA
ACGAA\$CGAG
GACGAA\$CGA
AGACGAA\$CG
GAGACGAA\$C

Permutations triées

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$ CGAGACGAA

A \$ CGAGACGA

AA \$ CGAGACG

GAA \$ CGAGAC

CGAA \$ CGAGA

ACGAA \$ CGAG

GACGAA \$ CGA

AGACGAA \$ CG

GAGACGAA \$ C

Permutations triées

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

A\$CGAGACGA

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

A\$CGAGACGA

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

GAA\$CGAGAC

CGAA\$CGAGA

ACGAA\$CGAG

GACGAA\$CGA

AGACGAA\$CG

GAGACGAA\$C

Permutations triées

\$CGAGACGAA

A\$CGAGACGA

AA\$CGAGACG

ACGAA\$CGAG

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$
\$CGAGACGAA
A\$CGAGACGA
AA\$CGAGACG
GAA\$CGAGAC
CGAA\$CGAGA
ACGAA\$CGAG
GACGAA\$CGA
AGACGAA\$CG
GAGACGAA\$C

Permutations triées

\$ CGAGACGAA
A \$ CGAGACGA
AA \$ CGAGACG
ACGAA \$ CGAG
AGACGAA \$ CG
CGAA \$ CGAGA
CGAGACGAA \$
GAA \$ CGAGAC
GACGAA \$ CGA
GAGACGAA \$ C

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA \$
\$ CGAGACGAA
A \$ CGAGACGA
AA \$ CGAGACG
GAA \$ CGAGAC
CGAA \$ CGAGA
ACGAA \$ CGAG
GACGAA \$ CGA
AGACGAA \$ CG
GAGACGAA \$ C

Permutations triées

\$ C G A G A C G A A
A \$ C G A G A C G A
A A \$ C G A G A C G
A C G A A \$ C G A G
A G A C G A A \$ C G
C G A A \$ C G A G A
C G A G A C G A A \$
G A A \$ C G A G A C
G A C G A A \$ C G A
G A G A C G A A \$ C

Transformée de Burrows-Wheeler

$T = \text{CGAGACGAA\$}$

Permutations circulaires

CGAGACGAA\$
\$CGAGACGAA
A\$CGAGACGA
AA\$CGAGACG
GAA\$CGAGAC
CGAA\$CGAGA
ACGAA\$CGAG
GACGAA\$CGA
AGACGAA\$CG
GAGACGAA\$C

Permutations triées

\$ CGAGACGA A
A \$ CGAGACG A
A A \$ CGAGAC G
ACGAA \$ CGA G
AGACGAA \$ CG G
CGAA \$ CGAG A
CGAGACGAA \$
GAA \$ CGAGAC C
GACGAA \$ CGA
GAGACGAA \$ C

$\text{TBW}(T) = \text{AAGGG\$CAC}$

Fonction LF

0 1 2 3 4 5 6 7 8 9
 $T = \text{CGAGACGAA\$}$

F	L
\$ C G A G A C G A A	
A \$ C G A G A C G A	
A A \$ C G A G A C G	
A C G A A \$ C G A G	
A G A C G A A \$ C G	
C G A A \$ C G A G A	
C G A G A C G A A \$	
G A A \$ C G A G A C	
G A C G A A \$ C G A	
G A G A C G A A \$ C	

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
\$	A
A	A
A	G
A	G
A	G
C	A
C	\$
G	C
G	A
G	C

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
$\$_1$	A_4
A_4	A_3
A_3	G_3
A_2	G_2
A_1	G_1
C_2	A_2
C_1	$\$_1$
G_3	C_2
G_2	A_1
G_1	C_1

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L	
$\$_1$	A_4	
A_4	A_3	Propriété
A_3	G_3	Les lettres identiques sont dans le
A_2	G_2	même ordre dans F et dans L .
A_1	G_1	
C_2	A_2	
C_1	$\$_1$	
G_3	C_2	
G_2	A_1	
G_1	C_1	

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L	
\$ ₁	A ₄	Propriété
A ₄	A ₃	Les lettres identiques sont dans le même ordre dans F et dans L .
A ₃	G ₃	
A ₂	G ₂	
A ₁	G ₁	Fonction LF
C ₂	A ₂	À partir de L et F , passage d'une lettre dans T à la précédente.
C ₁	\$ ₁	
G ₃	C ₂	
G ₂	A ₁	
G ₁	C ₁	

Récupération du texte

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$1	A ₄
A ₄	A ₃
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$1	A4
A4	A3
A3	G3
A2	G2
A1	G1
C2	A2
C1	\$1
G3	C2
G2	A1
G1	C1

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$A

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$ ₁	A ₄
A ₄	A ₄
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$A

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$ ₁	A ₄
A ₄	A ₃
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$AA

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$ ₁	A ₄
A ₄	A ₃
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$AA

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
$\$_1$	A_4
A_4	A_3
A_3	G_3
A_2	G_2
A_1	G_1
C_2	A_2
C_1	$\$_1$
G_3	C_2
G_2	A_1
G_1	C_1

Propriété

Les lettres identiques sont dans le même ordre dans F et dans L .

Fonction LF

À partir de L et F , passage d'une lettre dans T à la précédente.

Récupération du texte

\$AAG

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
$\$_1$	A_4
A_4	A_3
A_3	G_3
A_2	G_2
A_1	G_1
C_2	A_2
C_1	$\$_1$
G_3	C_2
G_2	A_1
G_1	C_1

Propriété

Les lettres identiques sont dans le même ordre dans F et dans L .

Fonction LF

À partir de L et F , passage d'une lettre dans T à la précédente.

Récupération du texte

\$AAG

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
$\$_1$	A_4
A_4	A_3
A_3	G_3
A_2	G_2
A_1	G_1
C_2	A_2
C_1	$\$_1$
G_3	C_2
G_2	A_1
G_1	C_1

Propriété

Les lettres identiques sont dans le même ordre dans F et dans L .

Fonction LF

À partir de L et F , passage d'une lettre dans T à la précédente.

Récupération du texte

\$AAGC

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
\$ ₁	A ₄
A ₄	A ₃
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans F et dans L .

Fonction LF

À partir de L et F , passage d'une lettre dans T à la précédente.

Récupération du texte

\$AAGC

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L
$\$_1$	A_4
A_4	A_3
A_3	G_3
A_2	G_2
A_1	G_1
C_2	A_2
C_1	$\$_1$
G_3	C_2
G_2	A_1
G_1	C_1

Propriété

Les lettres identiques sont dans le même ordre dans F et dans L .

Fonction LF

À partir de L et F , passage d'une lettre dans T à la précédente.

Récupération du texte

\$AAGCA

Fonction *LF*

$T = \text{CGAGACGAA\$}$

<i>F</i>	<i>L</i>
\$ ₁	A ₄
A ₄	A ₃
A ₃	G ₃
A ₂	G ₂
A ₁	G ₁
C ₂	A ₂
C ₁	\$ ₁
G ₃	C ₂
G ₂	A ₁
G ₁	C ₁

Propriété

Les lettres identiques sont dans le même ordre dans *F* et dans *L*.

Fonction *LF*

À partir de *L* et *F*, passage d'une lettre dans *T* à la précédente.

Récupération du texte

\$AAGCA

Fonction LF

$T = \text{CGAGACGAA\$}$

F	L	
$\$_1$	A_4	Propriété
A_4	A_3	Les lettres identiques sont dans le même ordre dans F et dans L .
A_3	G_3	
A_2	G_2	
A_1	G_1	Fonction LF
C_2	A_2	À partir de L et F , passage d'une lettre dans T à la précédente.
C_1	$\$_1$	
G_3	C_2	
G_2	A_1	
G_1	C_1	

Récupération du texte

\$AAGCA ...

FM-index

- ▶ Introduit par Ferragina et Manzini (2000) ;

FM-index

- ▶ Introduit par Ferragina et Manzini (2000) ;
- ▶ structure d'indexation compressée basée sur la transformée de Burrows-Wheeler ;

FM-index

- ▶ Introduit par Ferragina et Manzini (2000) ;
- ▶ structure d'indexation compressée basée sur la transformée de Burrows-Wheeler ;
- ▶ utilise :
 - ▶ la transformée de Burrows-Wheeler (lettres) ;
 - ▶ un échantillon de la table des suffixes (positions).

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

F	L
\$	CGAGACGA A
A	\$ CGAGACGA
A	A \$ CGAGAC G
A	CGAA \$ CGA G
A	GACGAA \$ CG G
C	GA A \$ CGAGA A
C	GAGACGAA \$
G	AA \$ CGAGA C
G	ACGAA \$ CG A
G	AGACGAA \$ C

Recherche de motifs

Recherchons $P = \underline{\text{A} \text{G} \text{A}}$ dans $T = \text{C} \text{G} \text{A} \text{G} \text{A} \text{C} \text{G} \text{A} \text{A} \$$

	0	1	2		0	1	2	3	4	5	6	7	8	9
F			↑	L										
	\$	C	G	A	G	A	C	G	A	A				
	A		\$	C	G	A	G	A	C	G	A			
	A	A	\$	C	G	A	G	A	C	G				
	A	C	G	A	A	\$	C	G	A	G				
	A	G	A	C	G	A	A	\$	C	G				
	C	G	A	A	\$	C	G	A	G	A				
	C	G	A	G	A	C	G	A	A	\$				
	G	A	A	\$	C	G	A	G	A	C				
	G	A	C	G	A	A	\$	C	G	A				
	G	A	G	A	C	G	A	A	\$	C				

Recherche de motifs

Recherchons $P = \underline{\text{A} \text{G} \text{A}}$ dans $T = \text{C} \text{G} \text{A} \text{G} \text{A} \text{C} \text{G} \text{A} \text{A} \$$

	0	1	2	3	4	5	6	7	8	9
F			↑							
L										

\$ C G A G A C G A A
A \$ C G A G A C G A
A A \$ C G A G A C G
A C G A A \$ C G A G
A G A C G A A \$ C G

C G A A \$ C G A G A
C G A G A C G A A \$
G A A \$ C G A G A C
G A C G A A \$ C G A
G A G A C G A A \$ C

Recherche de motifs

Recherchons $P = \underline{A} \textcolor{red}{G} \underline{A}$ dans $T = \text{CGAGACGAA\$}$

	0	1	2	3	4	5	6	7	8	9
F	A									
L		C	G	A	G	A	C	G	A	A
\$	C	G	A	G	A	C	G	A	A	
A	\$	C	G	A	G	A	C	G	A	
A	A	\$	C	G	A	G	A	C	G	
A	C	G	A	A	\$	C	G	A	C	
A	G	A	C	G	A	A	\$	C	G	
C	G	A	A	\$	C	G	A	G	A	
C	G	A	G	A	C	G	A	A	\$	
G	A	A	\$	C	G	A	G	A	C	
G	A	C	G	A	A	\$	C	G	A	
G	A	G	A	C	G	A	A	\$	C	

Recherche de motifs

Recherchons $P = \underline{A \textcolor{red}{G} A}$ dans $T = \text{CGAGACGAA\$}$

F L

\$ C G A G A C G A A ↑ Zéro G

A	\$	C	G	A	G	A	C	G	A
A	A	\$	C	G	A	G	A	C	G
A	C	G	A	A	\$	C	G	A	G
A	G	A	C	G	A	A	\$	C	G

Trois G

C G A A \$ C G A G A

C G A G A C G A A \$

G A A \$ C G A G A C

G A C G A A \$ C G A

G A G A C G A A \$ C

Recherche de motifs

Recherchons $P = \underline{A} \textcolor{red}{G} \underline{A}$ dans $T = \text{CGAGACGAA\$}$

F	L
	0 1 2 3 4 5 6 7 8 9
\$	CGAGACGA A
A	\$ CGAGACG A
A	A \$ CGAGAC G
A	CGAA \$ CGA G
A	GACGAA \$ CG G
C	GA A \$ CGAG A
Zéro G	\uparrow C GAGACGAA \$
Trois G	{ G A A \$ CGAGA C G A C GAA \$ CG A G A G A C G A A \$ C }

Recherche de motifs

Recherchons $P = \underline{\text{A} \text{G} \text{A}}$ dans $T = \text{C} \text{G} \text{A} \text{G} \text{A} \text{C} \text{G} \text{A} \text{A} \$$

	0	1	2	3	4	5	6	7	8	9
F	A									
L										
	$\$$	C	G	A	G	A	C	G	A	A
		A	$\$$	C	G	A	G	A	C	G
		A	A	$\$$	C	G	A	G	A	G
		A	C	G	A	A	$\$$	C	G	G
		A	G	A	C	G	A	$\$$	C	G
		C	G	A	A	$\$$	C	G	A	A
		C	G	A	G	A	C	G	A	$\$$
		G	A	A	$\$$	C	G	A	G	C
		G	A	C	G	A	$\$$	C	G	A
		G	A	G	A	C	G	A	$\$$	C

Recherche de motifs

Recherchons $P = \underline{\text{A} \text{G} \text{A}}$ dans $T = \text{C} \text{G} \text{A} \text{G} \text{A} \text{C} \text{G} \text{A} \text{A} \$$

	0	1	2	3	4	5	6	7	8	9
F	A									
L										
	$\$$	C	G	A	G	A	C	G	A	A
		A	$\$$	C	G	A	G	A	C	G
		A	A	$\$$	C	G	A	G	A	G
		A	C	G	A	A	$\$$	C	G	G
		A	G	A	C	G	A	$\$$	C	G
		C	G	A	A	$\$$	C	G	A	A
		C	G	A	G	A	C	G	A	A
		C	G	A	G	A	C	G	A	A
		G	A	A	$\$$	C	G	A	G	A
		G	A	G	A	$\$$	C	G	A	A
		G	A	G	A	C	G	A	$\$$	C

\uparrow Trois A

\uparrow Un A

Recherche de motifs

Recherchons $P = \begin{matrix} 0 & 1 & 2 \\ \textcolor{red}{A} & \textcolor{red}{G} & \textcolor{red}{A} \end{matrix}$ dans $T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \textcolor{gray}{C} & \textcolor{gray}{G} & \textcolor{gray}{A} & \textcolor{gray}{G} & \textcolor{gray}{A} & \textcolor{gray}{C} & \textcolor{gray}{G} & \textcolor{gray}{A} & \textcolor{gray}{A} & \$ \end{matrix}$

\$ C G A G A C G A A
A \$ C G A G A C G A
A A \$ C G A G A C G
Trois A↑A C G A A \$ C G A G
Un A { **A G A C G A A \$ C G**
C G A A \$ C G A G A
C G A G A C G A A \$
G A A \$ C G A G A C
G A C G A A \$ C G A
G A G A C G A A \$ C

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

F	L
\$	CGAGACGA A
A	\$ CGAGACG A
A	A \$ CGAGAC G
A	CGAA \$ CGA G
A	G A CGAA \$ C G
C	GA A \$ CGAG A
C	GAGACGAA \$
G	AA \$ CGAGA C
G	ACGAA \$ CGA
G	AGACGAA \$ C

Il y a une seule occurrence de AGA dans T .

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

	F	L
9	\$	C G A G A C G A A
8	A	\$ C G A G A C G A
7	A	A \$ C G A G A C G
4	A	C G A A \$ C G A G
2	A	G A C G A A \$ C G
5	C	G A A \$ C G A G A
0	C	G A G A C G A A \$
6	G	A A \$ C G A G A C
3	G	A C G A A \$ C G A
1	G	A G A C G A A \$ C

Il y a une seule occurrence de AGA dans T .

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

	F	L
9	\$	CGAGACGA A
8	A	\$ CGAGACG A
7	A	A \$ CGAGAC G
4	A	CGAA \$ CGA G
2	A	G A CGAA \$ C G
5	C	GA A \$ CGAG A
0	C	GA GACGAA \$
6	G	AA \$ CGAGA C
3	G	ACGAA \$ CGA
1	G	AGACGAA \$ C

Il y a une seule occurrence de AGA dans T .

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

F L

9	\$ C G A G A C G A A
8	A \$ C G A G A C G A
7	A A \$ C G A G A C G
4	A C G A A \$ C G A G
2	A G A C G A A \$ C G
5	C G A A \$ C G A G A
0	C G A G A C G A A \$
6	G A A \$ C G A G A C
3	G A C G A A \$ C G A
1	G A G A C G A A \$ C

Table des suffixes

Il y a une seule occurrence de AGA dans T .

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

F L

9	\$ C G A G A C G A A
8	A \$ C G A G A C G A
7	A A \$ C G A G A C G
4	A C G A A \$ C G A G
2	A G A C G A A \$ C G
5	C G A A \$ C G A G A
0	C G A G A C G A A \$
6	G A A \$ C G A G A C
3	G A C G A A \$ C G A
1	G A G A C G A A \$ C

Table des suffixes ←



Trop d'espace



Échantillonnage

Il y a une seule occurrence de AGA dans T .

Recherche de motifs

Recherchons $P = \text{AGA}$ dans $T = \text{CGAGACGAA\$}$

	F	L
9	\$ C G A G A C G A A	A
	A \$ C G A G A C G A	
	A A \$ C G A G A C G	
	A C G A A \$ C G A G	
	A G A C G A A \$ C G	
	C G A A \$ C G A G A	
0	C G A G A C G A A \$	
6	G A A \$ C G A G A C	
3	G A C G A A \$ C G A	
	G A G A C G A A \$ C	

Il y a une seule occurrence de AGA dans T .