

Aligning NGS reads and SNP finding



Philippe Bardou & Christophe Klopp

Organisation

Morning (9h30 -12h30) :

- Sequence quality
 - Theory + exercises
- Read mapping
 - Theory + exercises

Afternoon (14h-17h) :

- SAM format
 - Theory + exercises
- Visualisation
 - Theory + exercises
- SNP calling
 - Theory + exercises

2

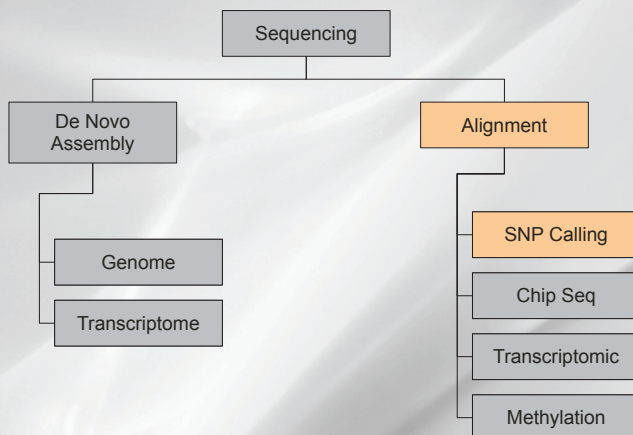
Why are you here?

Because

- You have produced or are going to produce a lot of reads for which you want to find variation between conditions or individuals or populations?
- You use :
 - Whole Genome Sequencing,
 - Sequence capture,
 - Reduced Representation Libraries,
 - Long range Polymerase Chain Reaction,
 - Amplification Fragment Length Polymorphisms.
- You use Solexa or Roche 454 platforms !
- You are curious!

3

Where are we?



What are you going to learn?

- To extract reads and reference genome from the NCBI
- To verify the read quality
- To format reference sequence
- To align the reads on the reference genome
- To index the reference sequence and the aligned reads
- To visualise the alignments and variations
- To call SNPs

What you should already know?

- How to connect to a remote unix server (putty)?
- What a unix command looks like?
- How to move around the unix environment?
- How to edit a file?

```

# ubuntu@localhost:~$ ssh ubuntu@192.168.1.100
Warning: Permanently added '192.168.1.100' (RSA) to the list of known hosts.
ubuntu@192.168.1.100:~$ cat /etc/passwd
root:x:0:0:root:/:/bin/bash
daemon:x:1:1:daemon:/usr/sbin:/usr/sbin/nologin
bin:x:2:2:bin:/usr/sbin:/usr/sbin/nologin
sys:x:3:3:sys:/dev:/usr/sbin/nologin
ubuntu:x:1000:1000:ubuntu:/home/ubuntu:/bin/bash
ubuntu@192.168.1.100:~$ cd /tmp
ubuntu@192.168.1.100:~/tmp$ touch file.txt
ubuntu@192.168.1.100:~/tmp$ ls -la
total 4
-rw-rw-r-- 1 ubuntu ubuntu 0 Oct 10 10:10 file.txt
ubuntu@192.168.1.100:~/tmp$ cat file.txt
ubuntu@192.168.1.100:~/tmp$ echo "hello" > file.txt
ubuntu@192.168.1.100:~/tmp$ cat file.txt
hello
ubuntu@192.168.1.100:~/tmp$ exit
ubuntu@192.168.1.100:~$
  
```

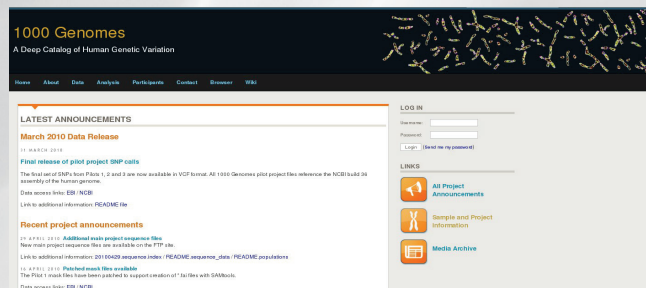
The pieces of software

- Quality : fastqc
- BWA : alignment
- Samtools : formatting SNP discovery
- IGV : visualisation

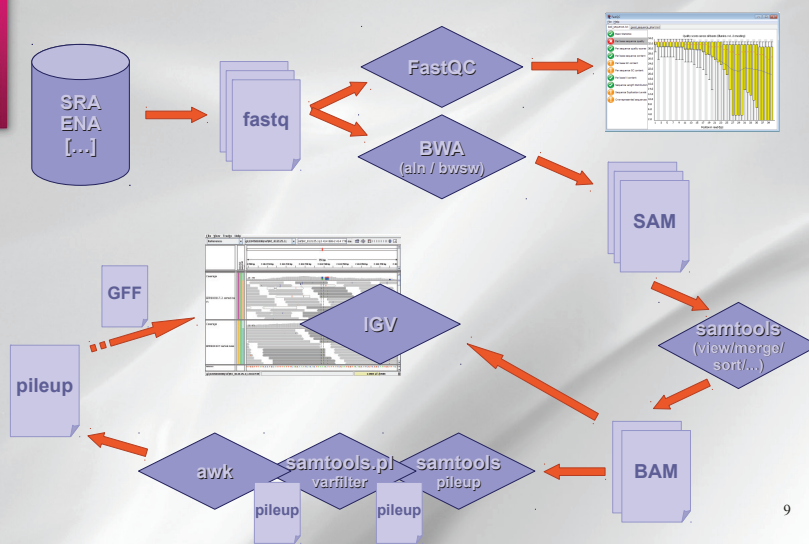


The 1000 genomes project

- Joint project NCBI / EBI
- Common data formats :
 - fastq
 - SAM (Sequence Alignment/Map)

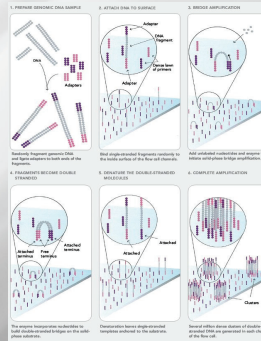
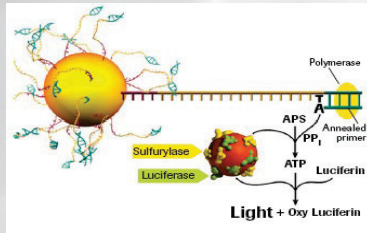


Overview



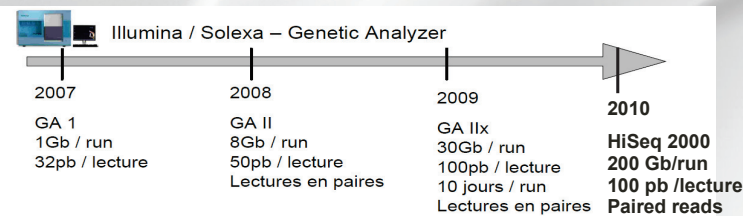
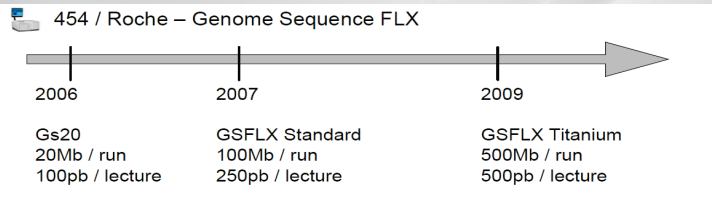
NGS platforms

- Two platforms :
 - Illumina Solexa
 - Roche 454



10

NGS reads



11

Sequencing bias bibliography

Research Open Access

Evaluation of next generation sequencing platforms for population targeted sequencing studies

Olivier Harismendy^{1,2*}, Pauline C Ng^{2,3*}, Robert L Strausberg², Xiaoyun Wang², Timothy B Stockwell², Karen Y Beeson¹, Nicholas J Schork², Sarah S Murray², Eric J Topol², Samuel Levy² and Kelly A Frazer^{2*}

Addresses: ¹ Scripps Genomic Medicine - Scripps Translational Science Institute - The Scripps Research Institute, N. Torrey Pines Court, La Jolla, CA 92037, USA. ² The J Craig Venter Institute, Medical Center Drive, Rockville, MD 20850, USA.

* These authors contributed equally to this work.

Correspondence: Samuel Levy. Email: slevy@jvci.org. Kelly A Frazer. Email: kfrazer@scripps.edu

Published: 27 March 2009
Genome Biology 2009, 10:R32 (doi:10.1186/gb-2009-10-3-r32)

Received: 14 December 2008
Revised: 23 February 2009
Accepted: 27 March 2009

The electronic version of this article is the complete one and can be found online at <http://genomebiology.com/2009/10/3/R32>

Published online 26 July 2008

Nucleic Acids Research, 2008, Vol. 36, No. 16 e105
doi:10.1093/nar/gkn425

Substantial biases in ultra-short read data sets from high-throughput DNA sequencing

Juliane C. Dohm¹, Claudio Lottaz², Tatiana Borodina¹ and Heinz Himmelbauer^{1,*}

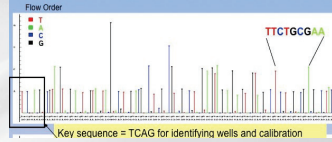
¹Max Planck Institute for Molecular Genetics, Ihnestr. 63-73, 14195 Berlin and ²Institute for Functional Genomics, Computational Diagnostics, University of Regensburg, Josef-Engert-Str. 9, 93053 Regensburg, Germany

Received December 21, 2007; Revised June 16, 2008; Accepted June 19, 2008

12

Sequencing bias

- Platform related
- Roche 454 (data from Jean-Marc Aury CNS)
 - 99,9% mapped reads
 - Mean error rate : 0,55%
 - 37% deletions, 53% insertions, 10% substitutions.
 - homopolymers errors
 - emPCR duplications
- Solexa (data from Jean-Marc Aury CNS)
 - 98,5% mapped reads
 - Mean error rate : 0,38%
 - 3% deletions, 2% insertions, 95% substitutions
 - Low A/T rich coverage

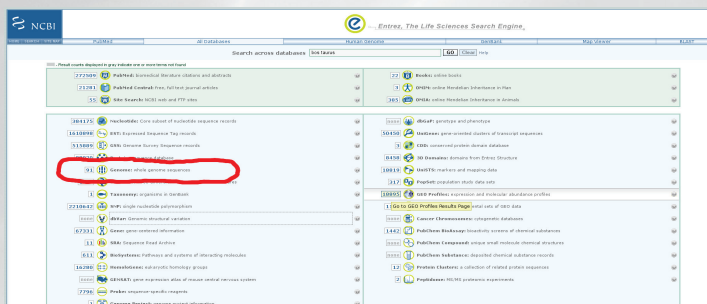


What data will we use?

- The needed data :
 - A reference sequence :
 - Genome
 - Parts of the genome
 - Transcriptome
 - Short reads

Where to get a reference genome?

- Assemble your own
- Use a public assembly :
 - NCBI : Genbank
 - EMBL



Where to get short reads?

- Produce your own sequences :
 - CNS
 - Local platform
 - Private company
- Use public data :
 - SRA : NCBI Sequence Read Archive
 - ENA : EMBL/EBI European Nucleotide Archive

NCBI SRA?

The screenshot shows the NCBI Entrez search engine interface. The search bar contains the text 'SRA'. Below the search bar, a list of search results is displayed. A red arrow points to the entry 'SRA: Sequence Read Archive', which is highlighted in blue. Other search results include 'PubMed: biomedical literature citations and abstracts', 'GenBank: genome survey sequence records', and 'SRA: Sequence Read Archive'.

EBI ENA

The screenshot shows the EBI ENA website. The page title is 'European Nucleotide Archive'. The main content area contains a search bar and a 'Text Search' section with a search query 'BN000065'. The search bar is labeled 'Enter search query, for example: BN000065' and has a 'Search' button. Below the search bar, there is a 'Sequence Search' section with a search query 'Enter or paste a nucleotide sequence' and a 'Search' button. The page also includes a 'Submit & update' section with a link to 'submission tools'.

- Meta data structure :
 - Experiment
 - Sample
 - Study
 - Run
 - Data file

ERP000014 Detecting variation in Salmonella Paratyphi A by sequencing pooled DNA

Accession	Spots	Bases
ERX000291	5.0M	193.7M
ERX000292	5.3M	191.0M
ERX000293	3.8M	129.6M
ERX000294	6.0M	201.9M
ERX000295	5.4M	195.4M
ERX000296	6.0M	215.9M
ERX000297	6.3M	209.4M

FASTQ format stores sequences and Phred qualities in a single file. It is concise and compact. FASTQ is first widely used in the Sanger Institute and therefore we usually take the Sanger specification and the standard FASTQ format, or simply FASTQ format. Although Solexa/Illumina read file looks pretty much like FASTQ, they are different in that the qualities are scaled differently. In the quality string, if you can see a character with its ASCII code higher than 90, probably your file is in the Solexa/Illumina format.

Example

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCTTCAGCGTTTCTCC
+
;;3;;;;;;;;;;7;;;;;;;;88
@EAS54_6_R1_2_1_540_792
TTGGCAGGCCAAGCCGATGGATCA
+
;;;;;;;;;;7;;;;-;3;83
@EAS54_6_R1_2_1_443_348
GTTGCTTCTGCCTGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
;;;;;;;;;;9;7;.7;393333
```

Published online 16 December 2009 *Nucleic Acids Research*, 2010, Vol. 38, No. 6 1767–1771
doi:10.1093/nar/gkp1137

SURVEY AND SUMMARY

The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants

Peter J. A. Cock^{1,*}, Christopher J. Fields², Naohisa Goto³, Michael L. Heuer⁴ and Peter M. Rice⁵

Table 1. The three described FASTQ variants, with columns giving the description, format name used in OBF projects, range of ASCII characters permitted in the quality string (in decimal notation), ASCII encoding offset, type of quality score encoded and the possible range of scores

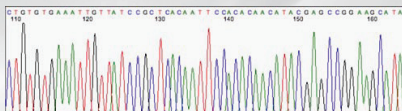
Description, OBF name	ASCII characters		Quality score	
	Range	Offset	Type	Range
Sanger standard				
fastq-sanger	33–126	33	PHRED	0 to 93
Solexa/early Illumina				
fastq-solexa	59–126	64	Solexa	-5 to 62
Illumina 1.3+				
fastq-illumina	64–126	64	PHRED	0 to 62

$$Q_{\text{PHRED}} = -10 \times \log_{10}(P_e)$$

$$Q_{\text{Solexa}} = -10 \times \log_{10}\left(\frac{P_e}{1 - P_e}\right)$$

Sequence quality

- Phred : base calling



What is Phred Quality?

Traditionally, Phred quality is defined on base calls. Each base call is an estimate of the true nucleotide. It is a random variable and can be wrong. The probability that a base call is wrong is called error probability.

Explanation about the quality values :

[source http://maq.sourceforge.net/qual.shtml](http://maq.sourceforge.net/qual.shtml)

22

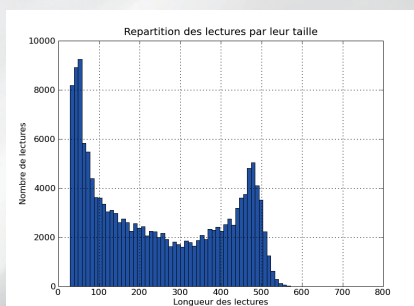
Which reads should I keep?

- All
- Some : what criteria and threshold should I use
 - Composition (number of Ns, complexity,...),
 - Quality,
 - Alignment based criteria,
- Should I trim the reads using :
 - Composition
 - Quality

23

Basic reads statistics

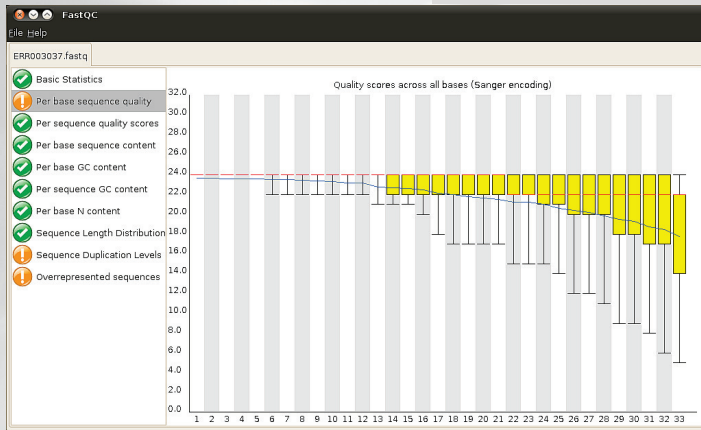
- Number of reads
- Length histogram
- Number of Ns in the reads
- Reads quality
- Reads redundancy
- Reads complexity



24

Sequence quality analysis

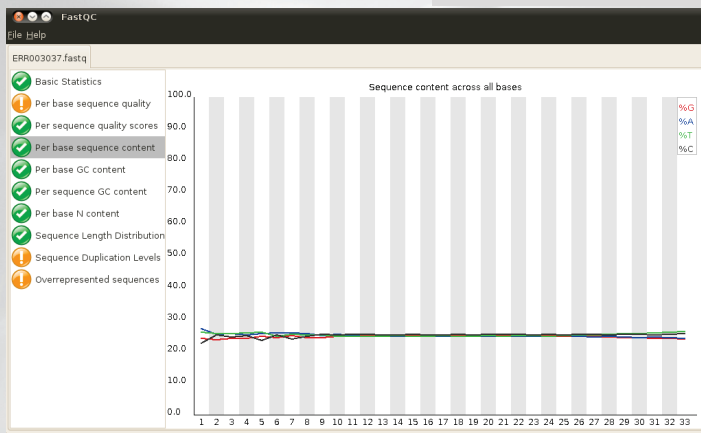
• FastQC :



25

Sequence quality analysis

• FastQC :



26

NG6 - Home

PROJETS RUNS
recherche OK

Présentation de NG6

NG6 est un environnement de stockage et de mise à disposition des données issues des nouvelles technologies des séquençage développé dans le cadre du partenariat liant les plate-formes bio-informatique et génomique Genotoul. Cet environnement, organisé autour des notions de projet et de run, permet un accès sécurisé aux données brutes, aux statistiques de traitements ainsi qu'aux assemblages et annotations produites. NG6 est un module typ3.

NG6 Presentation

NG6 is a storage platform aiming to provide an access to new sequencing generation data. This project has been initiated by the collaboration between the bioinformatics et genomic platform of Genotoul. This environment is organised around the notion of project and run, and allow a secured acces to raw data, statistics, assembly and produced annotations. NG6 is a typ3 modul.

Identification de l'utilisateur

Entrez votre nom d'utilisateur et votre mot de passe pour vous identifier:

Nom d'utilisateur:

Mot de passe:

Mot de passe oublié?

NG6 (Next Generation Sequencing Information System) développé par la plateforme Bioinformatique de Toulouse.

NG6 - Runs

Nom du run	Nom du projet	Date	Espèce	Nature des données	Type	Nombre de séquences	Taille totale des séquences	Description	Séquenceur
Phix validation	Démonstration	16-12-10	Phix	gDNA	1/8 ieme de flowcell A - lane 1	19091396	19282273196	4ieme run de validation, librairie illumina. Paired-end 2x100pb. Taille Insert = 250pb	HiSeq 2000
Démonstration - région 1	Démonstration	07-09-09	Escherichia Coli	ADNg	1/2 plaque	485812	100417483	Run public téléchargeable sur le site du NCBI	454 GS FLX
Démonstration - région 2	Démonstration	07-09-09	Escherichia Coli	ADNg	1/2 plaque	445381	104733228	Run public téléchargeable sur le site du NCBI	454 Titanium
E coli K12 région 1	Démonstration	28-05-09	Escherichia coli	ADNg	1/2 plaque	671856	291988221	Librairie Roche	454 Titanium
E coli K12 région 2	Démonstration	28-05-09	Escherichia coli	ADNg	1/2 plaque	529653	220797170	Librairie Roche	454 Titanium

28

NG6 - Run

Runs > Phix validation

Run Phix validation :
4ieme run de validation, librairie illumina. Paired-end 2x100pb. Taille Insert = 250pb

Nom du projet : Démonstration
 Date : 16-12-10
 Espèce : Phix
 Type : 1/8 ieme de flowcell A - lane 1
 Nature des données : gDNA
 Nombre de séquences : 19091396
 Taille totale des séquences : 19282273196
 Séquenceur : HiSeq 2000

Telechargements :
 s_1_phix.fasta.sorted.bam.bai
 s_1_phix.fasta.sorted.bam

Analyses réalisées :

Nom	Description	Logiciel	Version
AlignmentStats	Statistiques alignement contre genome de reference.	samtools flagstat	0.1.17 (r540)
ContaminationSearch	Recherche de contaminants sur les banques.	BWA	0.5.8c (r1536)
ReadsStats	Statistiques sur les lectures et leurs qualites.	fastqc	0.7.0

NG6 - Stats

Runs > Phix validation > ReadsStats

Analyse ReadsStats : Statistiques sur les lectures et leurs qualites.

Statistiques sur les lectures et sur leur qualité

Echantillons	Statistiques par position				Statistiques par séquence						
	Qualité	GCC	Ns	Contenue	Nombre de séquence	Qualité	GCC	Longueur	Niveau de duplication	Items	Séquences surreprésentées
s_1_1_seq(1)	PASS	PASS	PASS	PASS	95456798	PASS	45(WARN)	101(PASS)	FAIL	WARN	PASS
s_1_2_seq(2)	FAIL	PASS	PASS	PASS	95456798	PASS	44(WARN)	101(PASS)	FAIL	WARN	PASS

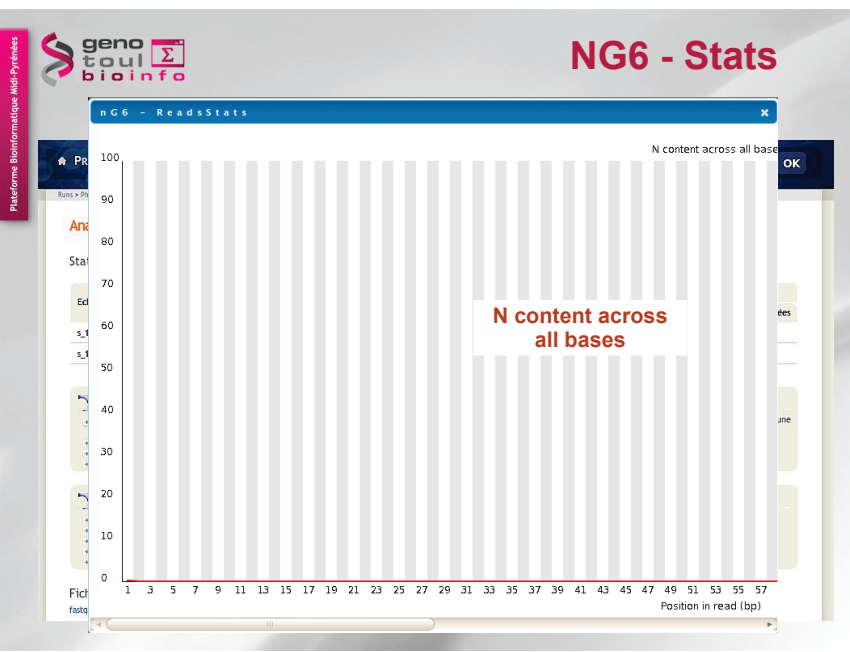
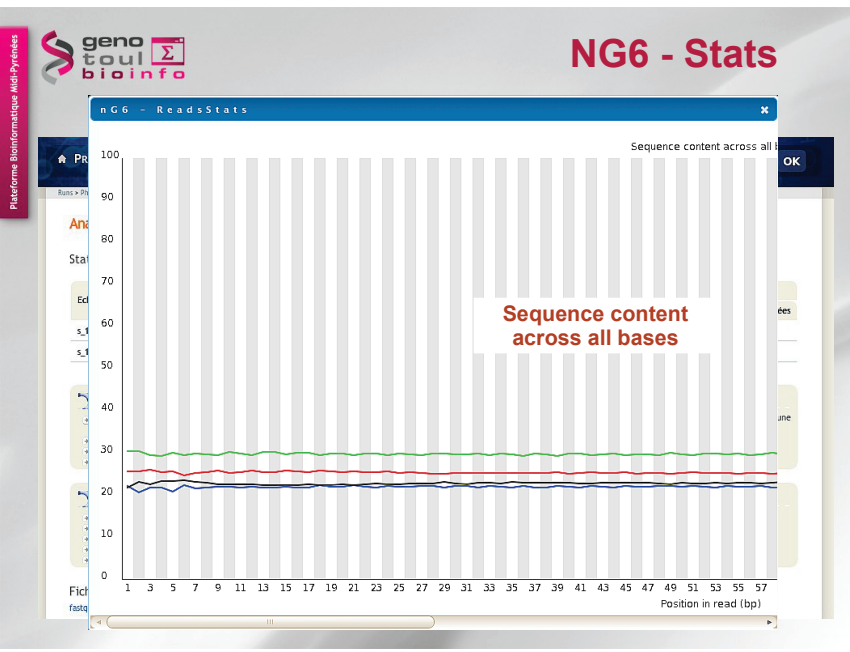
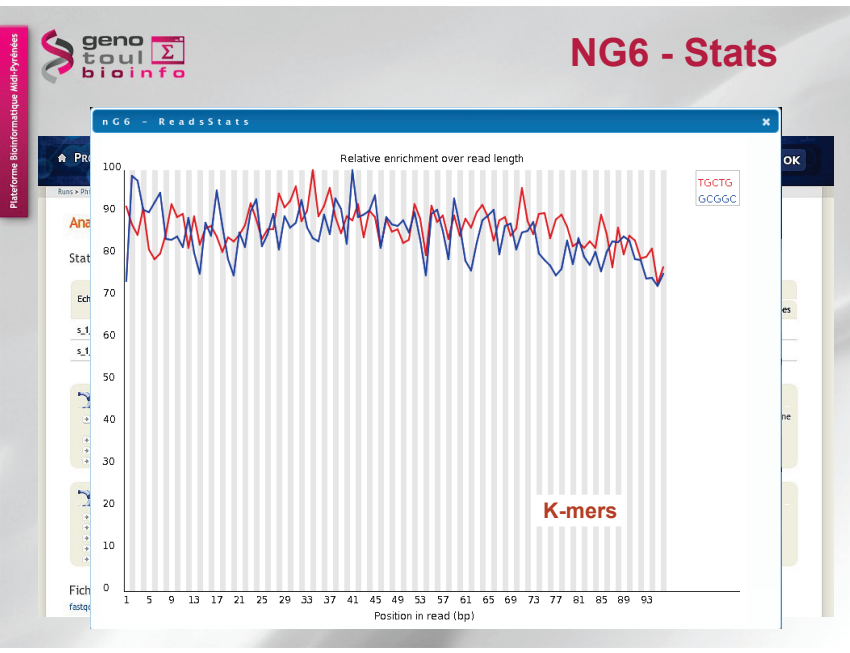
Aide Statistiques par position :

- Qualité : WARN = plus petit quartile d'une base inférieur à 10, ou si la mediane d'une base est inférieur à 25. FAIL = plus petit quartile d'une base inférieur à 5, ou si la mediane d'une base est inférieur à 20.
- GCC : WARN = GCC% par position +/- 5% de la moyenne. FAIL = GCC% par position +/- 10% de la moyenne.
- Ns : WARN = Ns% > 5%. FAIL = Ns% > 10%.
- Contenue : WARN = différence entre A et T ou G et C > 10% sur une position. FAIL = différence entre A et T ou G et C > 20% sur une position.

Aide Statistiques par séquence :

- Qualité : WARN = la qualité la plus observée - 27 (taux d'erreur de 0.2%). FAIL = la qualité la plus observée - 20 (taux d'erreur de 1%).
- GCC : WARN = plus de 15% des séquences ont un GCC différent de la distribution normale. FAIL = plus de 30% des séquences ont un GCC différent de la distribution normale.
- Longueur : WARN = toutes les séquences n'ont pas la même longueur. FAIL = une séquence a une longueur de 0pb.
- Niveau de duplication : WARN = plus de 20% des séquences sont non uniques. FAIL = plus de 30% des séquences sont non uniques.
- Séquences surreprésentées : WARN = une séquence représente plus de 0.1% du total. FAIL = une séquence représente plus de 1% du total.

Fichiers résultats
fastqc.tar.gz



Exercises / set 1

- snp server connexion
- Short read retrieval (wget)
- Read statistics (fastqc)

• Data sets :

- SRX002048
- ERR003037
- ERR000017

```

# Installation de FASTQ
# Ce script permet de télécharger les données de séquençage de la plateforme bioinformatique
# Vous disposez d'un quota initial de 100 Mo sur la partition /home (révisable de 0
# Vous disposez d'un espace de travail temporaire de 100 Mo partagé et non limité
# Le serveur "ftp" étant partagé par tous les utilisateurs, merci de ne pas y
# Informations concernant le cluster de calcul
# Utilisation du cluster de calcul (CC) se fait grâce à la commande "qsub" ou
# Vous disposez d'un espace de travail temporaire de 100 Mo partagé et non limité
# Règles de fonctionnement
# Mettez les points verra programme sur un petit jeu de données sur le serveur "
# Détails de l'usage de calcul en publiant vos fichiers de données sur /home
# Télécharger vos fichiers résultats (que vous souhaitez conserver) sur /home
Support
# Pour plus d'informations, consultez le site web :
# http://bioinfo.genotoul.fr
# Pour toute demande de support, contactez-vous à :
# support@genotoul.fr
# contact@genotoul.fr

```

Read alignment

- What are the ideas?
- The different software generations :
 - Smith-Waterman / Needleman-Wunch (1970)
 - BLAST (1990)
 - MAQ (2008)
 - BWA (2009)

BWA

- Much faster than MAQ
- Exact match
- Limited number of errors (2 for 32bp, 4 for 100 bp)

<http://bio-bwa.sourceforge.net/>

BIOINFORMATICS ORIGINAL PAPER Vol. 25 no. 14 2009, pages 1754–1760
doi:10.1093/bioinformatics/btp324

Sequence analysis

Fast and accurate short read alignment with Burrows–Wheeler transform

Heng Li and Richard Durbin*

Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK

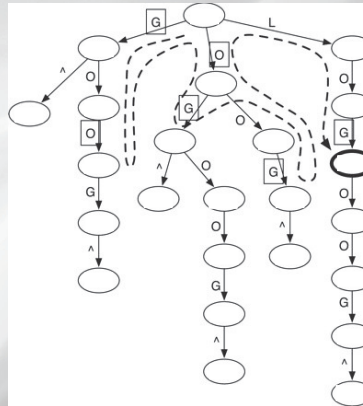
Received on February 20, 2009; revised on May 6, 2009; accepted on May 12, 2009

Advance Access publication May 18, 2009

Associate Editor: John Quackenbush

BWA prefix trie

- Word 'googol'
- ^ = start character
- --- Search of 'lol' with one error



The prefix trie is compressed to fit in memory in most cases (1Go for the human genome).

37

<http://bio-bwa.sourceforge.net/bwa.shtml>

Manual Reference Pages - bwa (1)

```

NAME
    bwa - Burrows-Wheeler Alignment Tool

CONTENTS
    Synopsis
    Description
    Commands And Options
    Sam Alignment Format
    Notes On Short-read Alignment
        Alignment Accuracy
        Estimating Insert Size Distribution
        Memory Requirement
        Speed
    Notes On Long-read Alignment
    See Also
    Author
    License And Citation
    History

SYNOPSIS
    bwa index -a bwtsv database.fasta
    bwa aln database.fasta short_read.fastq > aln_sa.sai
    bwa samse database.fasta aln_sa.sai short_read.fastq > aln.sam
    bwa sampe database.fasta aln_sa1.sai aln_sa2.sai read1.fq read2.fq > aln.sam
    bwa bwsw database.fasta long_read.fastq > aln.sam
    
```

38

Commands

Reference sequence indexing :

```
bwa index -a bwtsv db.fasta
```

Read Alignment :

```
bwa aln db.fasta short_read.fastq > aln_sa.sai
```

```
bwa bwsw database.fasta long_read.fastq > aln.sam
```

Formatting unpaired reads :

```
bwa samse db.fasta aln_sa.sai short_read.fastq > aln.sam
```

Formatting pair ends :

```
bwa sampe database.fasta aln_sa1.sai aln_sa2.sai read1.fq
read2.fq > aln.sam
```

39


```
index bwa index [-p prefix] [-a algoType] [-c] <in.db.fasta>
```

Index database sequences in the FASTA format.

OPTIONS:

- c Build color-space index. The input fast should be in nucleotide space.
- p STR Prefix of the output database [same as db filename]
- a STR Algorithm for constructing BWT index. Available options are:

is IS linear-time algorithm for constructing suffix array. It requires $5.37N$ memory where N is the size of the database. IS is moderately fast, but does not work with database larger than 2GB. IS is the default algorithm due to its simplicity. The current codes for IS algorithm are reimplemented by Yuta Mori.

bwtsw Algorithm implemented in BWT-SW. This method works with the whole human genome, but it does not work with database smaller than 10MB and it is usually slower than IS.

40

```
aln bwa aln [-n maxDiff] [-o maxGap] [-e maxGapE] [-d nbelTail] [-i nindelEnd] [-k maxSeedDiff] [-l seedLen] [-t nThreads] [-rni] [-M misMsc] [-o gapOsc] [-E gapEsc] [-q trinQual] <in.db.fasta> <in.query fq> > <out.sam>
```

Find the SA coordinates of the input reads. Maximum `maxSeedDiff` differences are allowed in the first `seedLen` subsequence and maximum `maxDiff` differences are allowed in the whole sequence.

OPTIONS:

- n NUM Maximum edit distance if the value is INT, or the fraction of missing alignments given 2% uniform base error rate if FLOAT. In the latter case, the maximum edit distance is automatically chosen for different read lengths. [0,04]
- o INT Maximum number of gap opens [1]
- e INT Maximum number of gap extensions, -1 for k-difference mode (disallowing long gaps) [-1]
- d INT Disallow a long deletion within INT bp towards the 3'-end [16]
- i INT Disallow an indel within INT bp towards the ends [5]
- l INT Take the first INT subsequence as seed. If INT is larger than the query sequence, seeding will be disabled. For long reads, this option is typically ranged from 25 to 35 for '-k 2'. [inf]
- k INT Maximum edit distance in the seed [2]
- t INT Number of threads (multi-threading mode) [1]
- M INT Mismatch penalty. BWA will not search for suboptimal hits with a score lower than (BestScore-misMsc). [3]
- O INT Gap open penalty [11]
- E INT Gap extension penalty [4]
- R INT Proceed with suboptimal alignments if there are no more than INT equally best hits. This option only affects paired-end mapping. Increasing this threshold helps to improve the pairing accuracy at the cost of speed, especially for short reads (-32bp).
- c Reverse query but not complement it, which is required for alignment in the color space.
- N Disable iterative search. All hits with no more than `maxDiff` differences will be found. This mode is much slower than the default.
- q INT Parameter for read trimming. BWA trims a read down to $\text{argmax}_x (\sum_{i=x+1}^l (\text{INT}-q_i))$ if $q_l < \text{INT}$ where l is the original read length. [0]

41

```
samse bwa samse [-n maxOcc] <in.db.fasta> <in.sam> <in.fq> > <out.sam>
```

Generate alignments in the SAM format given single-end reads. Repetitive hits will be randomly chosen.

OPTIONS:

- n INT Maximum number of alignments to output in the XA tag for reads paired properly. If a read has more than INT hits, the XA tag will not be written. [3]

```
sampe bwa sampe [-a maxInsSize] [-o maxOcc] [-n maxHitPaired] [-N maxHitDis] [-P] <in.db.fasta> <in1.sam> <in2.sam> <in1.fq> <in2.fq> > <out.sam>
```

Generate alignments in the SAM format given paired-end reads. Repetitive read pairs will be placed randomly.

OPTIONS:

- a INT Maximum insert size for a read pair to be considered being mapped properly. Since 0.4.5, this option is only used when there are not enough good alignment to infer the distribution of insert sizes. [500]
- o INT Maximum occurrences of a read for pairing. A read with more occurrences will be treated as a single-end read. Reducing this parameter helps faster pairing. [1000000]
- P Load the entire PH-index into memory to reduce disk operations (base-space reads only). With this option, at least 1.25N bytes of memory are required, where N is the length of the genome.
- n INT Maximum number of alignments to output in the XA tag for reads paired properly. If a read has more than INT hits, the XA tag will not be written. [3]
- N INT Maximum number of alignments to output in the XA tag for discordant read pairs (excluding singletons). If a read has more than INT hits, the XA tag will not be written. [10]

42

```

bwasw bwa bwasw [-a matchScore] [-b mmPen] [-q gapOpenPen] [-r gapExtPen] [-t nThreads]
[-w bandwidth] [-T thres] [-s hspIntv] [-z zBest] [-N nhspRev] [-c thresCoef]
<in.db.fasta> <in.fq>

```

Align query sequences in the <in.fq> file.

OPTIONS:

```

-a INT Score of a match [1]
-b INT Mismatch penalty [3]
-q INT Gap open penalty [5]
-r INT Gap extension penalty. The penalty for a contiguous gap of size k is
q*k*r. [2]
-t INT Number of threads in the multi-threading mode [1]
-w INT Band width in the banded alignment [33]
-T INT Minimum score threshold divided by a [37]
-c FLOAT Coefficient for threshold adjustment according to query length. Given an
l-long query, the threshold for a hit to be retained is
a*max(T,c*log(l)). [5.5]
-z INT Z-best heuristics. Higher -z increases accuracy at the cost of speed.
[1]
-s INT Maximum SA interval size for initiating a seed. Higher -s increases
accuracy at the cost of speed. [3]
-N INT Minimum number of seeds supporting the resultant alignment to skip
reverse alignment. [5]

```

43

Exercises / set 2

- Retrieving the reference sequence in fasta format :
 - `gj|224581838|ref|NC_012125.1|` Salmonella enterica subsp. enterica serovar Paratyphi C strain RKS4594, complete genome
- Indexing the reference sequence
- Aligning the reads (fastq format)
- Formatting the alignment in SAM

44

Sequence Alignment/Map (SAM) format

- Data sharing was a major issue with the 1000 genomes
- Capture all of the critical information about NGS data in a single indexed and compressed file
- Sharing : data across and tools
- Generic alignment format
- Supports short and long reads (454 – Solexa – Solid)
- Flexible in style, compact in size, efficient in random access

Website :

<http://samtools.sourceforge.net>

Paper :

Li H.*, Handsaker B.*, Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*, 25, 2078-9. [PMID: 19505943]

45

Sequence Alignment/Map (SAM) format

Aligners natively generating SAM

- **BFAST**, 'Blat-like Fast Accurate Search Tool' for Illumina and SOLiD reads.
- **Bowtie**, Highly efficient short read aligner. Natively support SAM output in recent version. A convertor is also available in samtools-C.
- **BWA**, Burrows-Wheeler Aligner for short and long reads.
- **GEM library**, Short read aligner. Convertor provided by the developers.
- **Karma**, the K-tuple Alignment with Rapid Matching Algorithm.
- **Mosaik**, The latest version support SAM output.
- **Novoalign**, An accurate aligner capable of gapped alignment for Illumina short reads. Academic free binary. Convertor is also available in samtools.
- **SNP-o-matic**, short read aligner and SNP caller.
- **SOLiD BaseQV Tool**, Developed by Applied Biosystems for converting SOLiD output files.
- **SSAHA2** (since v2.4), Classical aligner for both short and long reads.
- **Stampy**, by **Gerton Lunter**, An accurate read aligner capable of gapped alignment for Illumina short reads. Used for indel discovery on the 1000 genomes data. Not released.
- **TopHat** for mapping short RNA-seq reads bridging exon junctions.

SAM format Header section

- Header lines start with @ followed by a two-letter TAG
- Header fields are TYPE:VALUE pairs

Type	Tag	Description	
HD - header	VN*	File format version.	
	SO	Sort order. Valid values are: <i>unsorted, queryname or coordinate</i> .	
	GO	Group order (full sorting is not imposed in a group). Valid values are: <i>none, query or reference</i> .	
SQ - Sequence dictionary	SN*	Sequence name. Unique among all sequence records in the file. The value of this field is used in alignment records.	
	LN*	Sequence length.	
	AS	Genome assembly identifier. Refers to the reference genome assembly in an unambiguous form. Example: HG18.	
	M5	MDS checksum of the sequence in the uppercase (gaps and space are removed)	
	UR	URI of the sequence	
	SP	Species.	
RG - read group	ID*	Unique read group identifier. The value of the ID field is used in the RG tags of alignment records.	
	SM*	Sample (use pool name where a pool is being sequenced)	
	LB	Library	
	DS	Description	
	PU	Platform unit (e.g. lane for Illumina or slide for SOLiD); should be a full, unambiguous identifier	
	PI	Predicted median insert size (maybe different from the actual median insert size)	
	CC	Name of sequencing center producing the read.	
	DT	Date the run was produced (ISO 8601 date or datetime).	
	PL	Platform/technology used to produce the read.	
			@HD VN:1.0
PG - Program	PN	Program name	@SQ SN:chr20 LN:62435964
	RV	Program version	@RG ID:L1 PU:SC_1_10 LB:SC_1 SM:NA12891
	CS	Command line	@RG ID:L2 PU:SC_2_12 LB:SC_2 SM:NA12891
CO - comment		One-line text comments	

SAM format Alignment section

- 11 mandatory fields
- Variable number of optional fields
- Fields are tab delimited

1. **QNAME**: Query name of the read or the read pair
2. **FLAG**: Bitwise flag (pairing, strand, mate strand, etc.)
3. **RNAME**: Reference sequence name
4. **POS**: 1-Based leftmost position of clipped alignment
5. **MAPQ**: Mapping quality (Phred-scaled)
6. **CIGAR**: Extended CIGAR string (operations: MIDNSHP)
7. **MRNM**: Mate reference name ('=' if same as RNAME)
8. **MPOS**: 1-based leftmost mate position
9. **ISIZE**: Inferred insert size
10. **SEQQuery**: Sequence on the same strand as the reference
11. **QUAL**: Query quality (ASCII-33=Phred base quality)

SAM format Full example

Header

```

ERR000017_2.sam
@SQ
SN:ref LN:4833080
1 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
2 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT <<<<<<<<<<<<<<< XT:A:R NM:i:2 MD:Z:5A5A6
3 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
4 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
5 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
6 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
7 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
8 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
9 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>>>>>>>>>>>> XT:A:R NM:i:2 MD:Z:5A5A6
10 0 ref 4702037 25 18M * 0 0 CTATGAGCTATATGTTT >>>>777>>777>< XT:A:U NM:i:2 MD:Z:3C11G2
11 16 ref 2919865 37 18M * 0 0 5GGTGTATGTCTTTC >>>>>>>>>>>>>>>> XT:A:U NM:i:0 MD:Z:18
12 0 ref 2996664 37 18M * 0 0 GTTTTGTATGTGATAT ;'>70>>7>3877*4(7 XT:A:U NM:i:0 MD:Z:18
13 16 ref 510805 37 18M * 0 0 ATTCTCTATGAGTGAGT </>+798+>>>>>>>> XT:A:U NM:i:0 MD:Z:18
14 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>>>>>>>>>>>> XT:A:R NM:i:2 MD:Z:5A5A6
15 4 + 0 0 * * 0 0 GTGACACTCTGCTCTG >8>1>>8-9.7/(458 XT:A:R NM:i:2 MD:Z:5A5A6
16 16 ref 740202 0 18M * 0 0 TTTTTTTTTTTTTTTTTT >>>>????????????? XT:A:R NM:i:2 MD:Z:5A5A6
17 0 ref 1847349 37 SM118M * 0 0 CATGCAATATATCAT >49;:8,77;+*6+/' XT:A:U NM:i:2 MD:Z:5T11
    
```

Alignment

<QNAME> <FLAG> <RNAME> <POS> <MAPQ> <CIGAR> <MRNM> <MPOS> <ISIZE> <SEQ> <QUAL>

[<TAG>:<VTYPE>:<VALUE> [...]]

X? : Reserved for end users
 NM : Number of nuc. Difference
 MD : String for mismatching positions
 RG : Read group
 [...]

A : Printable character
 i : Signed 32-bit integer
 f : Single-precision float number
 Z : Printable string
 H : Hex string (high nybble first)

SAM format Flag field

Flag	Description
0x0001	the read is paired in sequencing, no matter whether it is mapped in a pair
0x0002	the read is mapped in a proper pair (depends on the protocol, normally inferred during alignment) ¹
0x0004	the query sequence itself is unmapped
0x0008	the mate is unmapped ¹
0x0010	strand of the query (0 for forward, 1 for reverse strand)
0x0020	strand of the mate ¹
0x0040	the read is the first read in a pair ^{1,2}
0x0080	the read is the second read in a pair ^{1,2}
0x0100	the alignment is not primary (a read having split hits may have multiple primary alignment records)
0x0200	the read fails platform/vendor quality checks
0x0400	the read is either a PCR duplicate or an optical duplicate

<http://picard.sourceforge.net/explain-flags.html>

SAM format Extended CIGAR format

Ref: GCATTAGATGCAGTACGC
 Read: ccTCAG--GCATTAgTg
 POS CIGAR
 5 2S4M2D6M3S

op	Description
M	Alignment match (can be a sequence match or mismatch)
I	Insertion to the reference
D	Deletion from the reference
N	Skipped region from the reference
s	Soft clip on the read (clipped sequence present in <seq>)
H	Hard clip on the read (clipped sequence NOT present in <seq>)
P	Padding (silent deletion from the padded reference sequence)

SAM format Extended CIGAR format

P	Padding (silent deletion from the padded reference sequence)
---	--

REF: CACGATCA**GACCGATACGTCCGA	REF: CACGATCA**GACCGATACGTCCGA
READ1: CGATCAGAGACCGATA	READ1: CGATCAGAGACCGATA
READ2: ATCA*AGACCGATAC	READ2: ATCAA*GACCGATAC
READ3: GATCA**GACCG	READ3: GATCA**GACCG
READ1: 6M2I8M	READ1: 6M2I8M
READ2: 4M1P1I9M	READ2: 4M1I1P9M
READ3: 5M2P5M	READ3: 5M2P5M

N	Skipped region from the reference
---	-----------------------------------

```
REF: AGCTAGCATCGTGTGCGCCGCTAGCATACGCATGATCGACTGTCAGCTAGTCAGACTAGTCGATCGATGTG
READ: GTGTAACCC.....TCAGAATA
```

where '.' on the read sequence indicates the intron. The CIGAR for this alignment is: 9M32N8M.

BAM format

- > Binary representation of SAM
- > Compressed by BGZF library
- > Greatly reduces storage space requirements to about 27% of original SAM

SAMtools

- > Library and software package
- > Creating sorted and indexed BAM files from SAM files
- > Removing PCR duplicates
- > Merging alignments
- > Visualization of alignments from BAM files
- > SNP calling
- > Short indel detection

<http://samtools.sourceforge.net/samtools.shtml>

SAMtools Example usage

```
inf212:~/Bureau/FormationNGS/JeuxDonnées> ../samtools-0.1.7a/samtools
Program: samtools (Tools for alignments in the SAM format)
Version: 0.1.7 (r510)

Usage: samtools <command> [options]

Command: view      SAM<->BAM conversion
          sort      sort alignment file
          pileup    generate pileup output
          faidx     index/extract FASTA
          tview     text alignment viewer
          index     index alignment
          fixmate   fix mate information
          glview    print GLPV3 file
          flagstat  simple stats
          calmd     recalculate MD/NM tags and
          merge     merge sorted alignments
          rmdup     remove PCR duplicates

inf212:~/Bureau/FormationNGS/JeuxDonnées> ../samtools-0.1.7a/samtools view
Usage: samtools view [options] <in.bam>[<in.sam>] [region1 [...]]

Options: -b      output BAM
         -h      print header for the SAM output
         -H      print header only (no alignments)
         -S      input is SAM
         -u      uncompressed BAM output (force -b)
         -x      output FLAG in HEX (samtools-C specific)
         -X      output FLAG in string (samtools-C specific)
         -t FILE list of reference names and lengths (force -S) [null]
         -T FILE reference sequence file (force -S) [null]
         -o FILE output file name [stdout]
         -f INT  required flag, 0 for unset [0]
         -F INT  filtering flag, 0 for unset [0]
         -q INT  minimum mapping quality [0]
         -L STR  only output reads in library STR [null]
         -r STR  only output reads in read group STR [null]
         -?      longer help
```

55

SAMtools Example usage

- > Create BAM from SAM


```
samtools view -bS aln.sam -o aln.bam
```
- > Sort BAM file


```
samtools sort example.bam sortedExample
```
- > Merge sorted BAM files


```
samtools merge sortedMerge.bam sorted1.bam sorted2.bam
```
- > Index BAM file


```
samtools index sortedExample.bam
```
- > Visualize BAM file


```
samtools tview sortedExample.bam reference.fa
```

56

Picard

- > A SAMtools complementary package
- > More format conversion than SAMtools
- > Visualization of alignments not available
- > SNP calling & short indel detection not available

<http://picard.sourceforge.net/>

57

Picard Example usage

- > ValidateSamFile
- > SortSam
- > MarkDuplicates
- > EstimateLibraryComplexity
- > MergeSamFiles
- > ViewSam
- > ReplaceSamHeader
- > SamToFastq
- > FastqToSam
- > SamFormatConverter
- > CreateSequenceDictionary
- > CleanSam
- > CompareSAMs

```
java -Xmx2g -jar PicardCmd.jar OPTION1=value1 OPTION2=value2...
```

Exercise / set 3

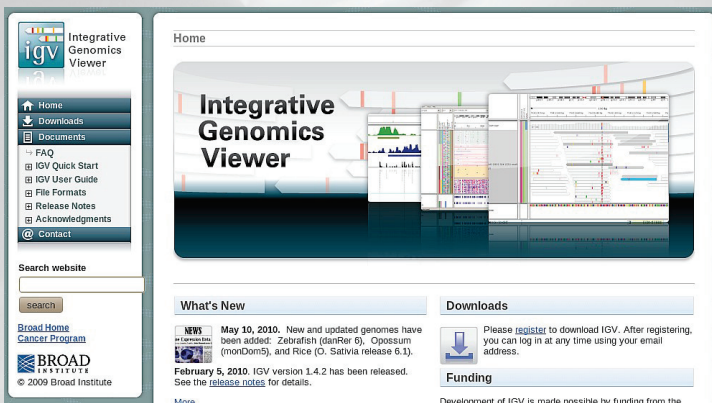
Visualizing the alignment SAMtools : tview

```
samtools tview aln.sorted.bam ref.fasta
```

The screenshot shows the terminal output of the command `samtools tview aln.sorted.bam ref.fasta`. The interface includes a menu with options like `g` (Go to specific location), `m` (Color for mapping quality), `n` (Color for nucleotide), `b` (Color for base quality), `c` (Color for cs qual), `z` (Color for ss qual), `o` (Toggle on/off dot view), `r` (Toggle on/off ref skip), `N` (Turn on at view), `C` (Turn on cs view), `i` (Toggle on/off ins), and `q` (Exit). A sample alignment record is displayed below the menu, showing a sequence of 'G' and 'A' characters and their corresponding quality scores.

Visualizing the alignment IGV

- IGV : Integrative Genomics Viewer
- Website : <http://www.broadinstitute.org/igv>



61

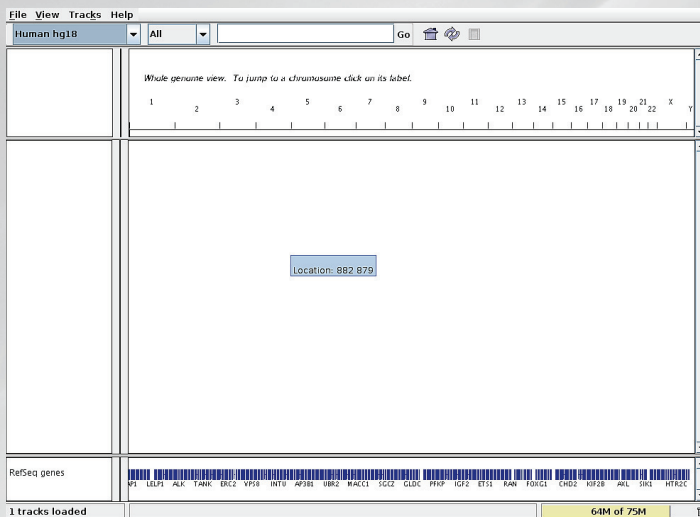
Visualizing the alignment IGV

- High-performance visualization tool
- Interactive exploration of large, integrated datasets
- Supports a wide variety of data types
- Documentations
- Developed at the Broad Institute of MIT and Harvard

File Formats
File Extension Identifies Format
Recommended File Formats
BAM
BED
CRS
CN
Cytoband
FASTA
GCT
genePred
GFF
gistic
HDP5
IGV
LOH
Readmate Files
MUT
RES
SAM
Sample Information
SEG
SNP
TAB
TDF
Track Line
Type Line
WIG

62

Visualizing the alignment IGV



63

Platforme Bioinformatique Midi-Pyrénées

geno
toul
bioinfo

Visualizing the alignment IGV - Loading the reference

File View Tracks Help

Load from File...
Load from URL...
Load from Server...
New Session...
Open Session...
Save Session...
Import Genome...
Remove Imported Genomes...
Save Image...
Export Regions...
Import Regions...
Clear Regions...
Exit
/home/bardou/igv_session.xml

Name

Sequence File * ...

Cytoband File ...

Gene File ...

* Required

The sequence file (required) can be a FASTA file, a directory of FASTA files, or a zip of FASTA files. Optionally, specify a cytoband file to display the chromosome ideogram and an annotation file to display the gene track. See the documentation for descriptions of supported annotation formats.

Save Cancel

RefSeq genes

1 tracks loaded 68M of 75M

64

Platforme Bioinformatique Midi-Pyrénées

geno
toul
bioinfo

Visualizing the alignment IGV - Loading the reference

File View Tracks Help

Reference g|224581838|ref|NC_012125.1| g|224581838|ref|NC_012125.1| Go

1 000 kb 2 000 kb 4 833 kb 3 000 kb 4 000 kb

Location: 1 814 674

Reference

1 tracks loaded 57M of 110M

65

Platforme Bioinformatique Midi-Pyrénées

geno
toul
bioinfo

Visualizing the alignment IGV - Loading the bam file

File View Tracks Help

Load from File... 24581838|ref|NC_012125.1| g|224581838|ref|NC_012125.1| Go

1 000 kb 2 000 kb 4 833 kb 3 000 kb 4 000 kb

Rechercher dans : JeuxDonnées

- Piqa_ERR000017_2
- Piqa_ERR000307
- ERR000017_2.bam
- ERR000017_2.fastq
- ERR000017_2.fastq.qualstat
- ERR000017_2.fastq.qualstat.boxplot.png
- ERR000017_2.fastq.qualstat.png
- ERR000017_2.sai
- ERR000017_2.sam
- ERR000017_2.sorted.bam
- ERR000017_2.sorted.bam.bai
- ERR000307.bam
- ERR000307.fastq
- ERR000307.fastq.qualstat
- ERR000307.fastq.qualstat.png
- ERR000307.sai
- ERR000307.sam
- ERR000307.sorted.bam
- ERR000307.sorted.bam.bai
- NC_012125.1.fasta
- NC_012125.1.gff

Nom de fichier : "ERR000017_2.sorted.bam" "ERR000307.sorted.bam"

Fichiers de type : Tous les fichiers

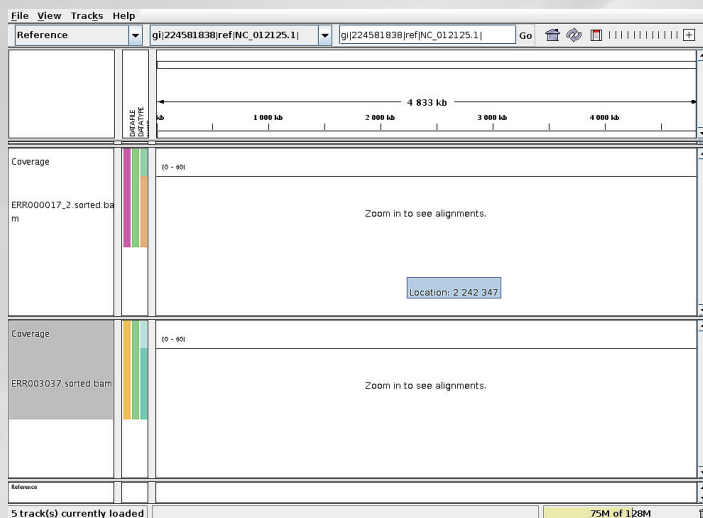
Ok Annuler

Reference

1 tracks loaded 59M of 110M

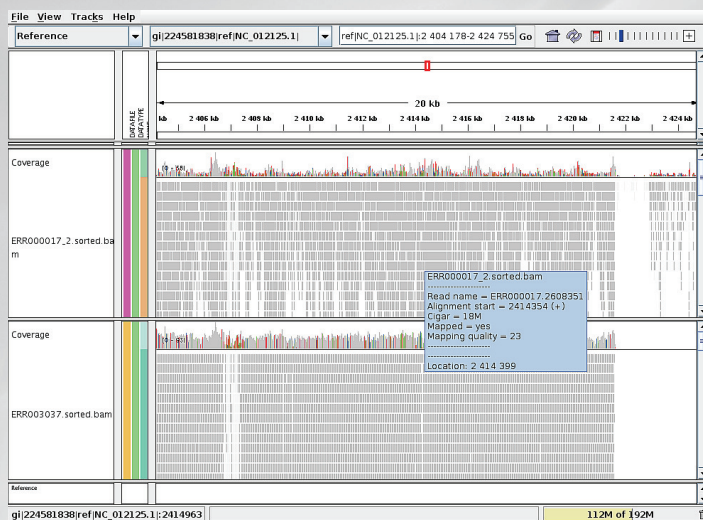
66

Visualizing the alignment IGV - Loading the bam file



67

Visualizing the alignment IGV - Zoom



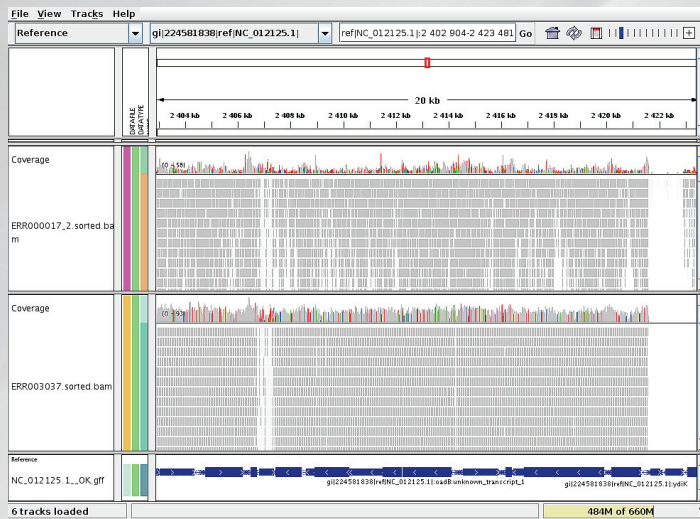
68

Visualizing the alignment IGV - Zoom



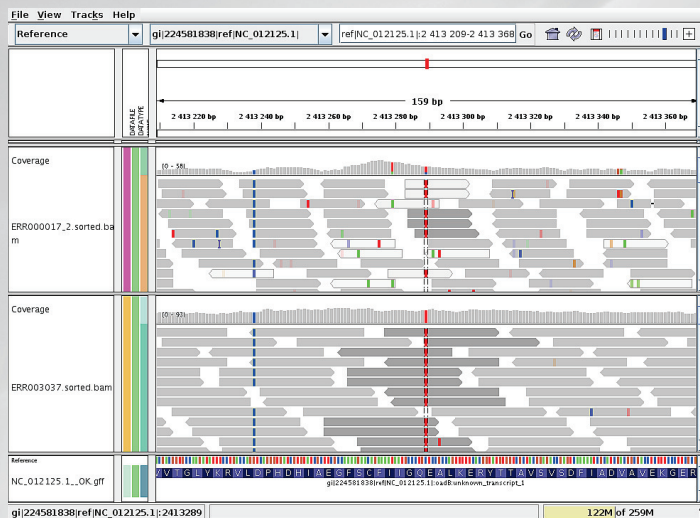
69

Visualizing the alignment IGV - Loading a gff file



70

Visualizing the alignment IGV - Loading a gff file



71

Visualizing the alignment IGV - Coverage

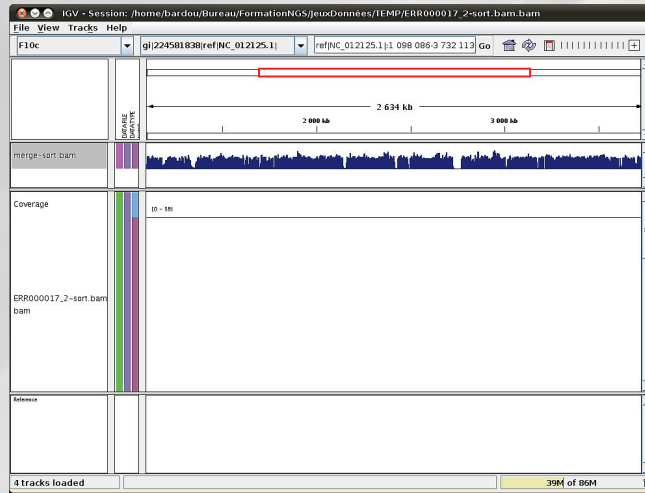
- > Generate the coverage information to be displayed in IGV.

```
java -jar igvtools.jar count aln.bam aln.depth.tdf ref.genome
```

- > Remark : ref.genome was generated when we imported the genome sequence
- > This step is optional, but it is essential if you want to see the read depth information in large scale.

72

Visualizing the alignment IGV - Coverage



Exercise / set 4

The pileup format

Chr. - Coord. - Base(** for indel) - Number of reads covering the site - **Read bases*** - Base qualities

```

seq1 272 T 24 .$. . . . . ^+. <<<+; <<<<<<<<<<=<; <7<&
seq1 273 T 23 .$. . . . . A <<<+; <<<<<<<<<3<=<<<+; <<+
seq1 274 T 23 .$. . . . . 7<7; <; <<<<<<<<=<; <; <<6
seq1 275 A 23 .$. . . . . ^\ . <+; 9* <<<<<<<<<<=<; <<<<<
seq1 276 G 22 .$. . . . . T. . . . . 33; +<<7=7<<7<<6<<1; <<6<
seq1 277 T 22 .$. . . . . C. . . . . G. . . . . +7<; <<<<<<<<<=<<+; <<6<
seq1 278 G 23 .$. . . . . ^k. %38* <<+; <7<<7<<=<<<+; <<<<<<
seq1 279 C 23 A..T. . . . . ; 75& <<<<<<<<<<<<<9<<<+; <<
seq2 156 A 11 .$. . . . . +2AG.+2AG.+2AGGG <975; ; <<<<<<

```

- Read bases :**
- > '.' and ';' : match to the reference base on the forward/reverse strand
 - > 'ACTGN' and 'actgn' : for a mismatch on the forward/reverse strand
 - > '^' and '\$' : start/end of a read segment
 - > '+[0-9]+[ACGTNacgtn]+' and '-[0-9]+[ACGTNacgtn]+' : insertion/deletion

<http://samtools.sourceforge.net/pileup.shtml>

Variant Calling with SAMtools

> Get the raw variant :

```
samtools pileup -vcf ref.fa aln.bam > raw.txt
samtools view -u aln.bam X | samtools pileup -vcf ref.fa - > raw-X.txt
```

- ref.fa Fasta formatted file of the reference genome
- aln.bam Sorted BAM formatted file, from the alignments
- raw[-X].txt Output pileup formatted, with consensus calls
- c Calls the consensus base at each position
- v Show positions that do not agree with ref.fa
- f Reference sequence, ref.fa (in FastA format)

The pileup format

Consensus base - Consensus quality - Probability of difference from ref. base - Max. mapping quality

seq1	60	T	T	66	0	99	13^~.^~.	9<<55<;<<<<<<
seq1	61	G	G	72	0	99	15^~.^y.	(;975&;<<<<<<<<
seq1	62	T	T	72	0	99	15	.\$.....	<;;55;<<<<<<<<
seq1	63	G	G	72	0	99	15	.\$.....^~.	4;2;<7;+<<<<<<<<
seq1	64	G	G	69	0	99	14	9+5<;;<<<<<<<<
seq1	65	A	A	69	0	99	14	.\$.....	<5-2<;<<<<<<<<
seq1	66	C	C	66	0	99	13	&*<;<<<<<<<8<
seq1	67	C	C	69	0	99	14^~.	,75<.4<<<<<<<<<
seq1	68	C	C	69	0	99	14	576<;7<<<<<8<<

seq2	156	A	A	10	0	99	11	.\$.....+2AG.+2AG.+2AGGG	<975;:<<<<<<
seq2	156	*	+AG/+AG	71	252	99	11	+AG * 3 8 0	

1st indel allele - 2nd indel allele - Reads supporting 1st - Reads supporting 2nd - Reads supporting 3rd
 Reads-bases Reads-qualities

Filter - samtools.pl

> Filter the raw variant calls :

```
samtools.pl varFilter raw.txt > raw_ok.txt
(samtools.pl varFilter -p raw.txt > raw_ok.txt) >& raw_filtered.txt
```

```
Usage: samtools.pl varFilter [options] <in.cns-pileup>
Options: -Q INT    minimum RMS mapping quality for SNPs [25]
         -q INT    minimum RMS mapping quality for gaps [10]
         -d INT    minimum read depth [3]
         -D INT    maximum read depth [100]
         -G INT    min indel score for nearby SNP filtering [25]
         -w INT    SNP within INT bp around a gap to be filtered [10]
         -W INT    window size for filtering dense SNPs [10]
         -N INT    max number of SNPs in a window [2]
         -l INT    window size for filtering adjacent gaps [30]
         -p       print filtered variants
```

Filter - awk

Acquire final variant calls by setting a quality threshold (50 for indels and 20 for substitutions) :

```
awk '($3=="*"&&$6>=50)||($3!="*"&&$6>=20)' raw.flt.txt > raw.final.txt
```

gi	ref	pos	mapq	flags	CIGAR	pos	RG	seq	qual		
gi_224581838	ref_NC_012125.1	58123	C	G	12	12	4	18	...GggggGggggggg	9AC4===CA?7c96e0=	
gi_224581838	ref_NC_012125.1	58137	G	G	50	0	10	20-lc,-lc,-lc,-lc,-lc,-lc,-lc,-lc,-lc,-lc	57C	
gi_224581838	ref_NC_012125.1	58137	*	C	/*	182	422	10	20	-C * 13 7 0 0 0	
gi_224581838	ref_NC_012125.1	58147	A	G	94	94	13	23	gGGGgggGggggggggggggg	<C?A?BcA&=<=c9?B99<e9	
gi_224581838	ref_NC_012125.1	58153	C	G	89	89	13	24	gGGGgggGggggggggggggg	CCA=A???==888<78067B?83	
gi_224581838	ref_NC_012125.1	58168	C	A	59	78	13	26	aaa.aaa.....a	@CB3AA>44@AAAGAc->?52@-D:	
gi_224581838	ref_NC_012125.1	58213	G	A	48	48	8	10	AAAAAaAaA	CAS@G?DOOO	
gi_224581838	ref_NC_012125.1	58222	A	G	21	21	8	10	.GGGGGGG,	@=8&A0;D	
gi_224581838	ref_NC_012125.1	58225	G	G	48	0	8	10	...+2AC.....	@?<7A83>0	
gi_224581838	ref_NC_012125.1	58225	*	*/AC	B	B	8	10	*+AC	B 1 0 0 0	
gi_224581838	ref_NC_012125.1	58231	G	A	21	21	8	10	.AAAAAA,,	C2, /<D9<0	
gi_224581838	ref_NC_012125.1	58235	G	C	49	49	7	12	CGGGCGGGGGCC	CABCA&A0D?Z	
gi_224581838	ref_NC_012125.1	58273	C	G	22	22	6	13	.GGGGG,,GgG	@5?A<BRADCC0CB	
gi_224581838	ref_NC_012125.1	58282	T	C	27	27	7	11	CCCC,,CCcC	CCTD0ADC1CC	

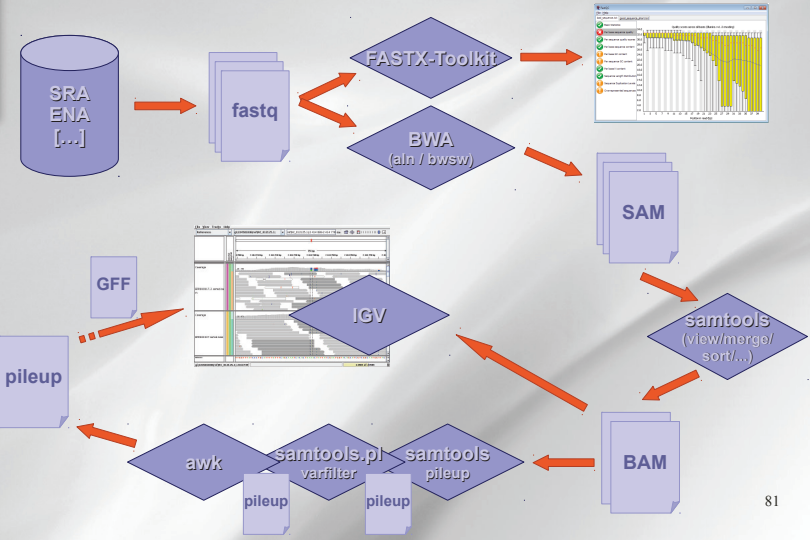
Consensus base - Consensus quality - Probability of difference from ref. base - Max. mapping quality

Consensus - samtools.pl

Consensus : A way of representing the results of a multiple sequence alignment (which residues are most abundant in the alignment at each position).

```
samtools.pl pileup2fq raw.txt > raw.fq
```

Synthesis



Exercise / set 5

82

The END

Satisfaction form :

<http://bioinfo.genotoul.fr/index.php?id=60>

Exam :

<http://bioinfo.genotoul.fr/index.php?id=122>

83