

Niniejsze opracowanie zostało stworzone przez dr Magdę Mielczarek, pracownika Uniwersytetu Przyrodniczego we Wrocławiu w ramach wykonywania obowiązków związanych z kształceniem studentów i jest przeznaczone dla studentów Biologii (Wydział Biologii i Hodowli Zwierząt) na potrzeby dydaktyczne bez prawa do dalszego rozpowszechniania.

OMÓWIENIE STANDARDOWYCH KROKÓW W ANALIZIE DANYCH NGS

Wykład 4

PORÓWNANIE METOD

	SANGER SEQ	NEXT GENERATION SEQ
STRATEGY	Separate reaction for the sequencing of all exons of a single gene / any piece of the genome	One single reaction for the simultaneous analysis of different genes / whole genomes!
ADVANTAGE	High reliability	Cost-effective and efficient by simultaneous and fast analysis Relatively cheap (0.001\$/1000 bp)
DISADVANTAGE	Expensive (0.1\$ / 1000 bp) Time consuming due to limited automation and necessity of many different reactions	Interpretation of the abundance of data is challenging High coverage needed for accuracy

KOSZT/CZAS SEKW. CAŁEGO GENOMU CZŁOWIEKA

- Projekt poznania genomu człowieka
 - ~ 13 lat
 - ~ 15 000 000 \$
- Technologia NGS
 - ~ kilka dni
 - ~ 2 000 \$

ANALIZA ZMIENNOŚCI GENETYCZNEJ ZA POMOCĄ WGS POPULARNE PLATFORMY

Illumina

- ~100-250 bp

Pacific Biosciences (PacBio)

- 10–25 kb

Oxford Nanopore Technologies

- 500 bp - 2.5 Mb

... i inne



GŁÓWNE ZASTOSOWANIA NGS

1. Detekcja mutacji i polimorfizmów genetycznych
2. Poznawanie nowych genomów (*de novo genome assembly*)
3. RNA-Seq: profilowanie transkryptomu
4. Chip-Seq: Interakcje na linii białko-DNA
5. Methyl-Seq: Epigenomika i metylacja DNA
6. Metagenomika

Xinkun Wang. *Next Generation Sequencing Data Analysis*. 2016, CRC PRESS

ANALIZA ZMIENNOŚCI GENETYCZNEJ ZA POMOCĄ WGS

Detection of genomic variation among individuals of a population is among the most frequent applications of next-generation sequencing (NGS).

Locating genomic sequence variations that correlate with disease predisposition or drug response, and establishing a genotypic basis of various phenotypes become common focuses of many NGS studies in biomedical and life sciences research.

Besides variations carried through the germline for generations, NGS has also been applied to identify de novo germline and somatic mutations, which occur more frequently than previously expected and underlie numerous human diseases including various types of cancer.

Xinkun Wang. *Next Generation Sequencing Data Analysis*. 2016, CRC PRESS

ANALIZA ZMIENNOŚCI GENETYCZNEJ ZA POMOCĄ WGS

Detecting the various forms of genetic variations/mutations from NGS (...) is not an easy task. The primary challenge is to differentiate true sequence variations/mutations from false positives caused by sequencing errors and artifacts generated in (...) sequence alignment.

It is, therefore, important to generate high-quality sequence data before performing data analysis. Equally important, sensitive and yet specific variant/mutant calling algorithms are required to achieve high accuracy in genomic variation and mutation discovery.

Xinkun Wang. *Next Generation Sequencing Data Analysis*. 2016, CRC PRESS

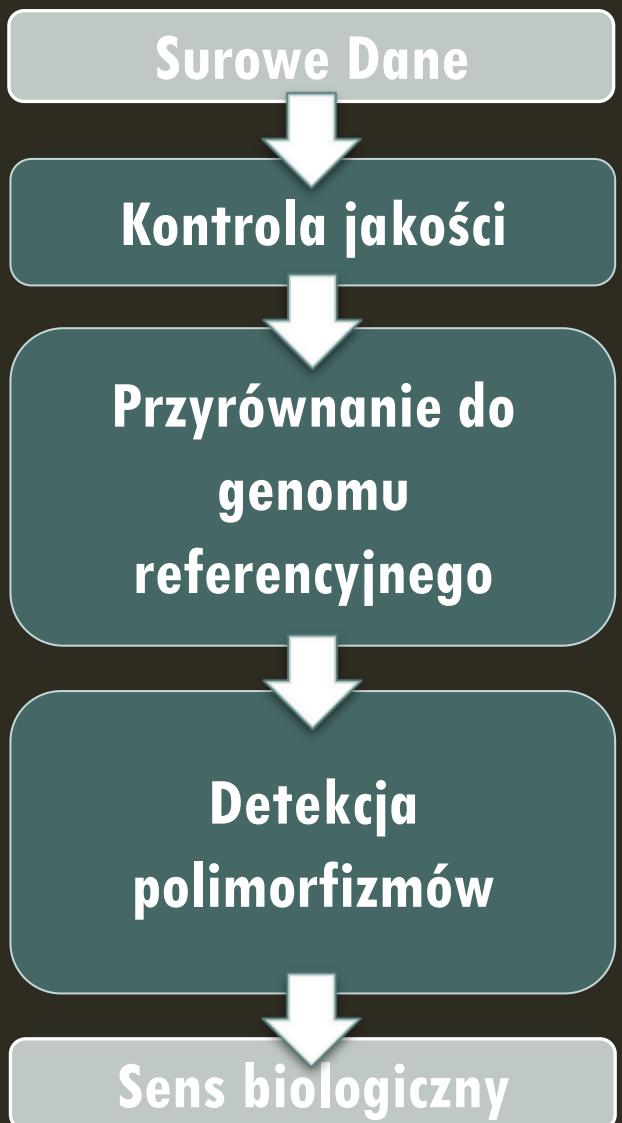
PIPELINE

Pipeline = łańcuch przetwarzania danych

Uproszczony schemat

Jedno z podstawowych zagadnień w analizie danych
NGS → przerównanie do genomu referencyjnego oraz
detekcja mutacji/polimorfizmów

Uproszczony schemat



SUROWE DANE

Read 1

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
AGAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAGGC  
+  
@ @@FFDFFGHGHFDDDGHHHDDDDHIIJJDDIIIGDDJGDDGD!  
(...)
```

Read 1

Single-end

AGAAATG...

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

SUROWE DANE

Read 1

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
AGAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAAGC  
+  
@ @@FFDFFGHGHFDDDGHHHDDDDHIIJJDDIIIGDDJGDDGD!  
(...)
```

Read 2

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 2:N:0:ACAGTG  
TTAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAAGC  
+  
DD@FF@ @FGHGHH01DDGHHHDDDDHIIJJJDIIIGDDJGDDGDD  
(...)
```



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

14G	Dec	6	14:35	wgs_S4505Nr1.1.fastq.gz
15G	Dec	6	14:56	wgs_S4505Nr1.2.fastq.gz
21G	Dec	6	15:27	wgs_S4505Nr10.1.fastq.gz
21G	Dec	6	15:59	wgs_S4505Nr10.2.fastq.gz
9.9G	Dec	6	16:14	wgs_S4505Nr11.1.fastq.gz
11G	Dec	6	16:30	wgs_S4505Nr11.2.fastq.gz
24G	Dec	6	17:06	wgs_S4505Nr12.1.fastq.gz
26G	Dec	6	17:45	wgs_S4505Nr12.2.fastq.gz
18G	Dec	6	18:12	wgs_S4505Nr13.1.fastq.gz
19G	Dec	6	18:40	wgs_S4505Nr13.2.fastq.gz
22G	Dec	6	19:12	wgs_S4505Nr14.1.fastq.gz
23G	Dec	6	19:46	wgs_S4505Nr14.2.fastq.gz
19G	Dec	6	20:14	wgs_S4505Nr15.1.fastq.gz
20G	Dec	6	20:44	wgs_S4505Nr15.2.fastq.gz
30G	Dec	6	21:28	wgs_S4505Nr16.1.fastq.gz
31G	Dec	6	22:15	wgs_S4505Nr16.2.fastq.gz
13G	Dec	6	22:34	wgs_S4505Nr17.1.fastq.gz
13G	Dec	6	22:53	wgs_S4505Nr17.2.fastq.gz
23G	Dec	6	23:27	wgs_S4505Nr18.1.fastq.gz
24G	Dec	7	00:02	wgs_S4505Nr18.2.fastq.gz
28G	Dec	7	00:44	wgs_S4505Nr19.1.fastq.gz
29G	Dec	7	01:28	wgs_S4505Nr19.2.fastq.gz
17G	Dec	7	01:53	wgs_S4505Nr2.1.fastq.gz
18G	Dec	7	02:20	wgs_S4505Nr2.2.fastq.gz
18G	Dec	7	02:46	wgs_S4505Nr20.1.fastq.gz
19G	Dec	7	03:14	wgs_S4505Nr20.2.fastq.gz
13G	Dec	7	03:33	wgs_S4505Nr3.1.fastq.gz
14G	Dec	7	03:54	wgs_S4505Nr3.2.fastq.gz
26G	Dec	7	04:33	wgs_S4505Nr4.1.fastq.gz
28G	Dec	7	05:15	wgs_S4505Nr4.2.fastq.gz

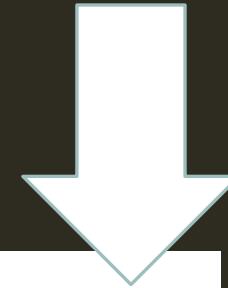
417M	Jan	16	2014	lane2_NoIndex_L002_R1_039.fastq.gz
395M	Jan	16	2014	lane2_NoIndex_L002_R1_044.fastq.gz
393M	Jan	16	2014	lane2_NoIndex_L002_R1_025.fastq.gz
404M	Jan	16	2014	lane2_NoIndex_L002_R1_008.fastq.gz
390M	Jan	16	2014	lane2_NoIndex_L002_R2_014.fastq.gz
391M	Jan	16	2014	lane2_NoIndex_L002_R2_046.fastq.gz
405M	Jan	16	2014	lane2_NoIndex_L002_R1_028.fastq.gz
395M	Jan	16	2014	lane2_NoIndex_L002_R1_024.fastq.gz
402M	Jan	16	2014	lane2_NoIndex_L002_R1_027.fastq.gz
406M	Jan	16	2014	lane2_NoIndex_L002_R1_058.fastq.gz
407M	Jan	16	2014	lane2_NoIndex_L002_R1_029.fastq.gz
392M	Jan	16	2014	lane2_NoIndex_L002_R1_015.fastq.gz
390M	Jan	16	2014	lane2_NoIndex_L002_R2_004.fastq.gz
392M	Jan	16	2014	lane2_NoIndex_L002_R1_052.fastq.gz
389M	Jan	16	2014	lane2_NoIndex_L002_R2_023.fastq.gz
404M	Jan	16	2014	lane2_NoIndex_L002_R1_040.fastq.gz
391M	Jan	16	2014	lane2_NoIndex_L002_R2_006.fastq.gz
394M	Jan	16	2014	lane2_NoIndex_L002_R1_054.fastq.gz
391M	Jan	16	2014	lane2_NoIndex_L002_R2_003.fastq.gz
399M	Jan	16	2014	lane2_NoIndex_L002_R1_006.fastq.gz
393M	Jan	16	2014	lane2_NoIndex_L002_R1_023.fastq.gz
384M	Jan	16	2014	lane2_NoIndex_L002_R2_045.fastq.gz
391M	Jan	16	2014	lane2_NoIndex_L002_R2_013.fastq.gz
410M	Jan	16	2014	lane2_NoIndex_L002_R1_030.fastq.gz
387M	Jan	16	2014	lane2_NoIndex_L002_R2_021.fastq.gz
397M	Jan	16	2014	lane2_NoIndex_L002_R1_003.fastq.gz

SUROWE DANE

PRZECHOWYWANIE DANYCH

```
total 44G
-rw-rw---- 1 szyda upwroclaw 383M Jan 16 2014 lanel_NoIndex_L001_R1_029.fastq.gz
-rw-rw---- 1 szyda upwroclaw 376M Jan 16 2014 lanel_NoIndex_L001_R2_055.fastq.gz
-rw-rw---- 1 szyda upwroclaw 385M Jan 16 2014 lanel_NoIndex_L001_R2_038.fastq.gz
-rw-rw---- 1 szyda upwroclaw 375M Jan 16 2014 lanel_NoIndex_L001_R1_046.fastq.gz
-rw-rw---- 1 szyda upwroclaw 387M Jan 16 2014 lanel_NoIndex_L001_R1_059.fastq.gz
-rw-rw---- 1 szyda upwroclaw 383M Jan 16 2014 lanel_NoIndex_L001_R2_048.fastq.gz
-rw-rw---- 1 szyda upwroclaw 385M Jan 16 2014 lanel_NoIndex_L001_R1_010.fastq.gz
-rw-rw---- 1 szyda upwroclaw 380M Jan 16 2014 lanel_NoIndex_L001_R2_050.fastq.gz
-rw-rw---- 1 szyda upwroclaw 378M Jan 16 2014 lanel_NoIndex_L001_R1_057.fastq.gz
```

1 genom



+ dane dodatkowe

@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG

AGAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAGGC

+

@@@FFDFFGHGHFDDGHHHHDDDDHIIJJDDIIIGDDJGDDGD!

(...)

@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>)
Y	Y if the read is filtered, N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence



Illumina

www.illumina.com

$$Q = -10 \log_{10} P \longrightarrow P = 10^{\frac{-Q}{10}}$$

KODOWANIE JAKOŚCI

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
AGAAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAGGC  
+  
@@@FFDFFGHGHHFDDDGHHHDDDDHIIJJDDIIIGDDJGDDGD!  
(...)
```

Phred Quality Score	Error	Accuracy (1 - Error)
10	1/10 = 10%	90%
20	1/100 = 1%	99%
30	1/1000 = 0.1%	99.9%
40	1/10000 = 0.01%	99.99%
50	1/100000 = 0.001%	99.999%
60	1/1000000 = 0.0001%	99.9999%



S - Sanger Phred+33, raw reads typically (0, 40)

X - Solexa Solexa+64, raw reads typically (-5, 40)

I - Illumina 1.3+: Phred+64, raw reads typically (0, 40)

J - Illumina 1.5+: Phred+64, raw reads typically (3, 41)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (**bold**)
(Note: See discussion above).

L - Illumina 1.8+: Phred+33, raw reads typically (0, 41)

N - Nanopore Phred+33, Duplex reads typically (0, 50)

E - ElemBio AVITI Phred+33, raw reads typically (0, 55)

P - PacBio Phred+33, HiFi reads typically (0, 93)

JAKOŚĆ NUKLEOTYDÓW

Read 1

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
AGAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAAGC  
+  
@@@FFDFFGHGHFDDDGHHHDDDDHIIJJDDIIIGDDJGDDGD!  
(...)
```

Read 2

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 2:N:0:ACAGTG  
TTAAATGCCAGGCTAGATGAGTTACAATCTAGTATCAAGATAAGC  
+  
DD@FF@@@FGHGHH01DDGHHHDDDDHIIJJJDIIIGDDJGDDGDD  
(...)
```

Surowe Dane

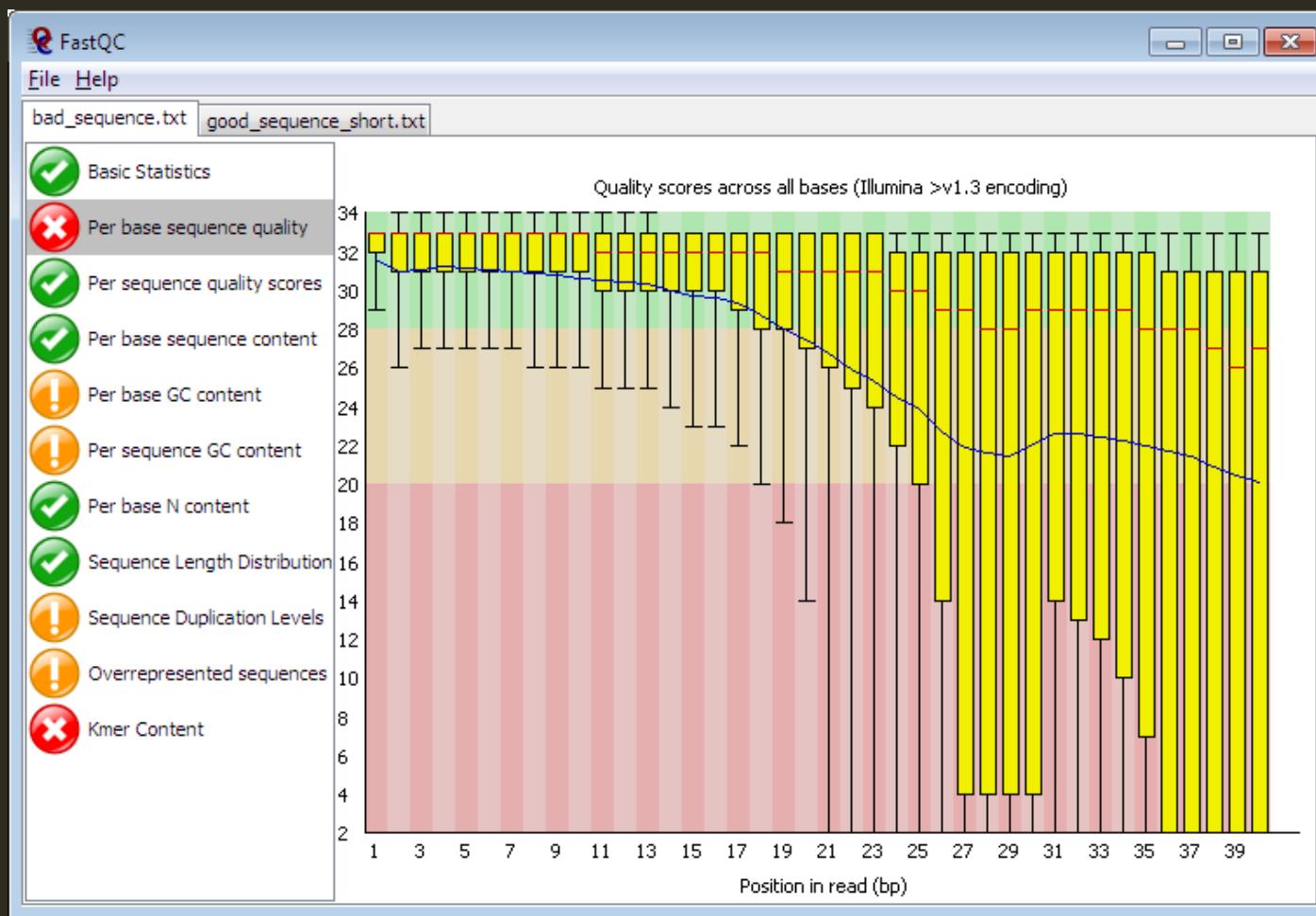
Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

KONTROLA JAKOŚCI



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

FASTQC



Babraham Bioinformatics

[About](#) | [People](#) | [Services](#) | [Projects](#) | [Training](#) | [Publications](#)

FastQC

Function

A quality control tool for high throughput sequence data.

The main functions of FastQC are

- Import of data from BAM, SAM or FastQ files (any variant)
- Providing a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive application

FASTQC - PRZYKŁADY

The main functions of FastQC are

- Import of data from BAM, SAM or FastQ files (any variant)
- Providing a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive appli



Documentation

A [copy of the FastQC documentation](#) is available for you to try before you buy (well download..).

Example Reports

- [Good Illumina Data](#)
- [Bad Illumina Data](#)
- [Adapter dimer contaminated run](#)
- [Small RNA with read-through adapter](#)
- [Reduced Representation BS-Seq](#)
- [PacBio](#)
- [454](#)

FASTQC - PRZYKŁADY

Summary

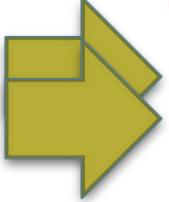
The main functions of FastQC are

- Import of data from BAM, SAM
- Providing a quick overview to t
- Summary graphs and tables to
- Export of results to an HTML b
- Offline operation to allow autor

Documentation

A [copy](#) of the FastQC documentation i

Example Reports



- [Good Illumina Data](#)
- [Bad Illumina Data](#)
- [Adapter dimer contaminated r](#)
- [Small RNA with read-through a](#)
- [Reduced Representation BS-S](#)
- [PacBio](#)
- [454](#)

Summary

- [Basic Statistics](#)
- [Per base sequence quality](#)
- [Per tile sequence quality](#)
- [Per sequence quality scores](#)
- [Per base sequence content](#)
- [Per sequence GC content](#)
- [Per base N content](#)
- [Sequence Length Distribution](#)
- [Sequence Duplication Levels](#)
- [Overrepresented sequences](#)
- [Adapter Content](#)

BASIC STATISTICS



Basic Statistics

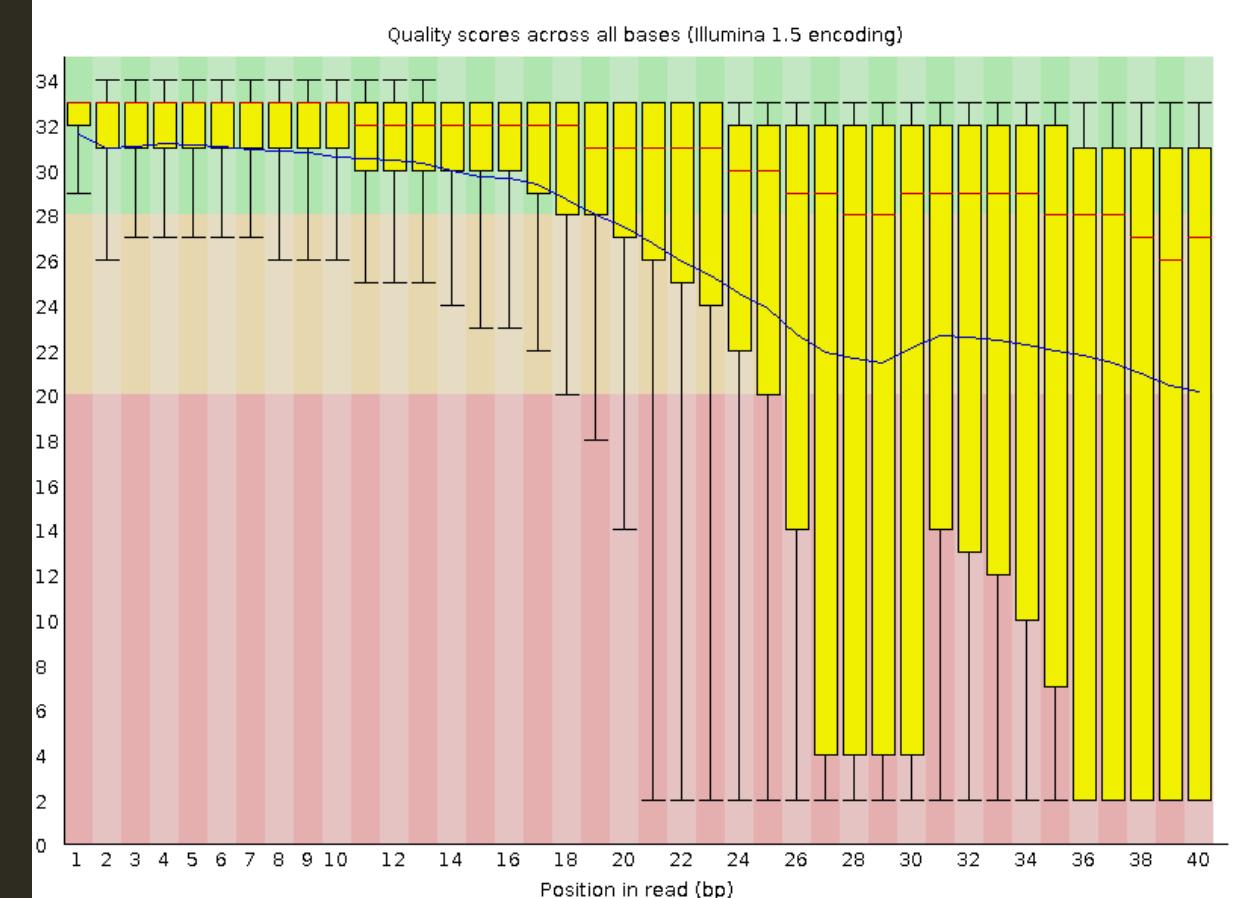
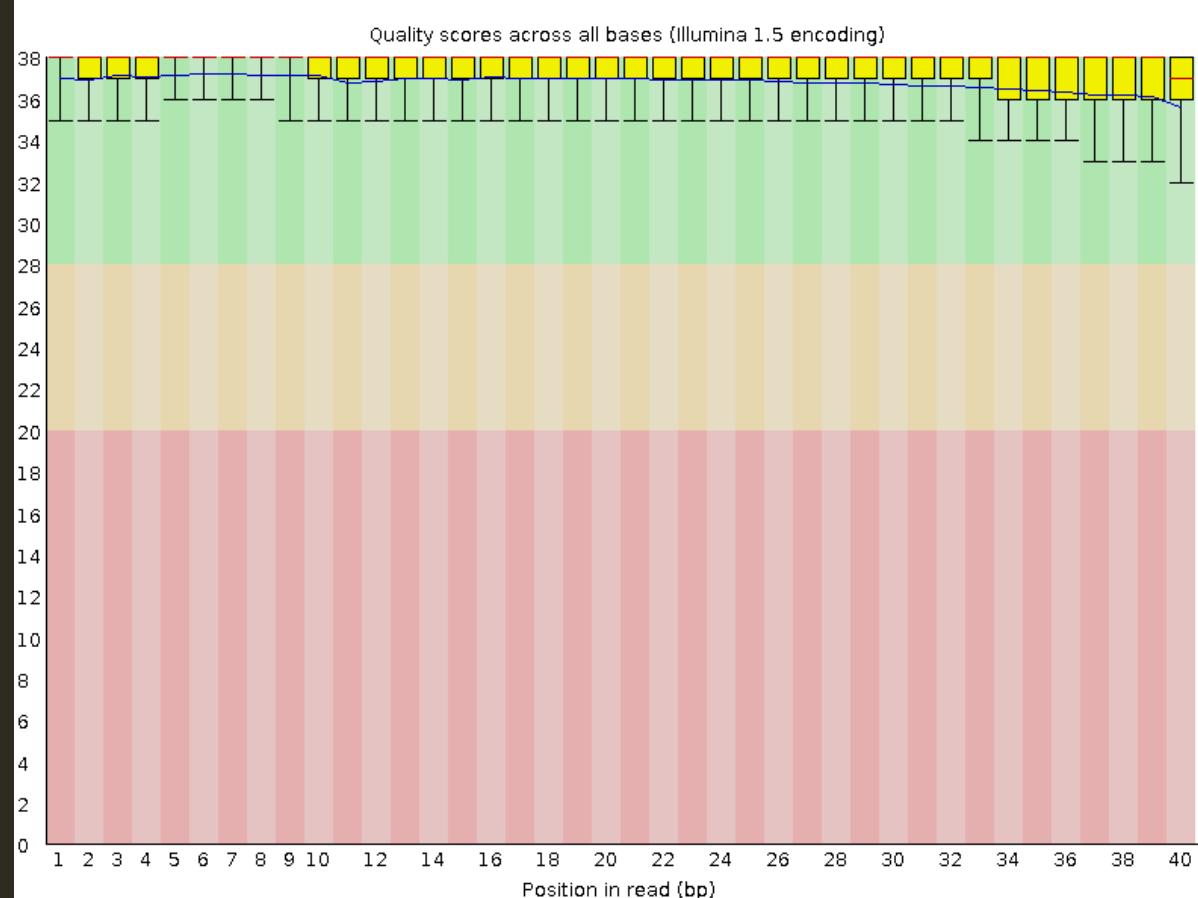
Measure	Value
Filename	good_sequence_short.txt
File type	Conventional base calls
Encoding	Illumina 1.5
Total Sequences	250000
Sequences flagged as poor quality	0
Sequence length	40
%GC	45



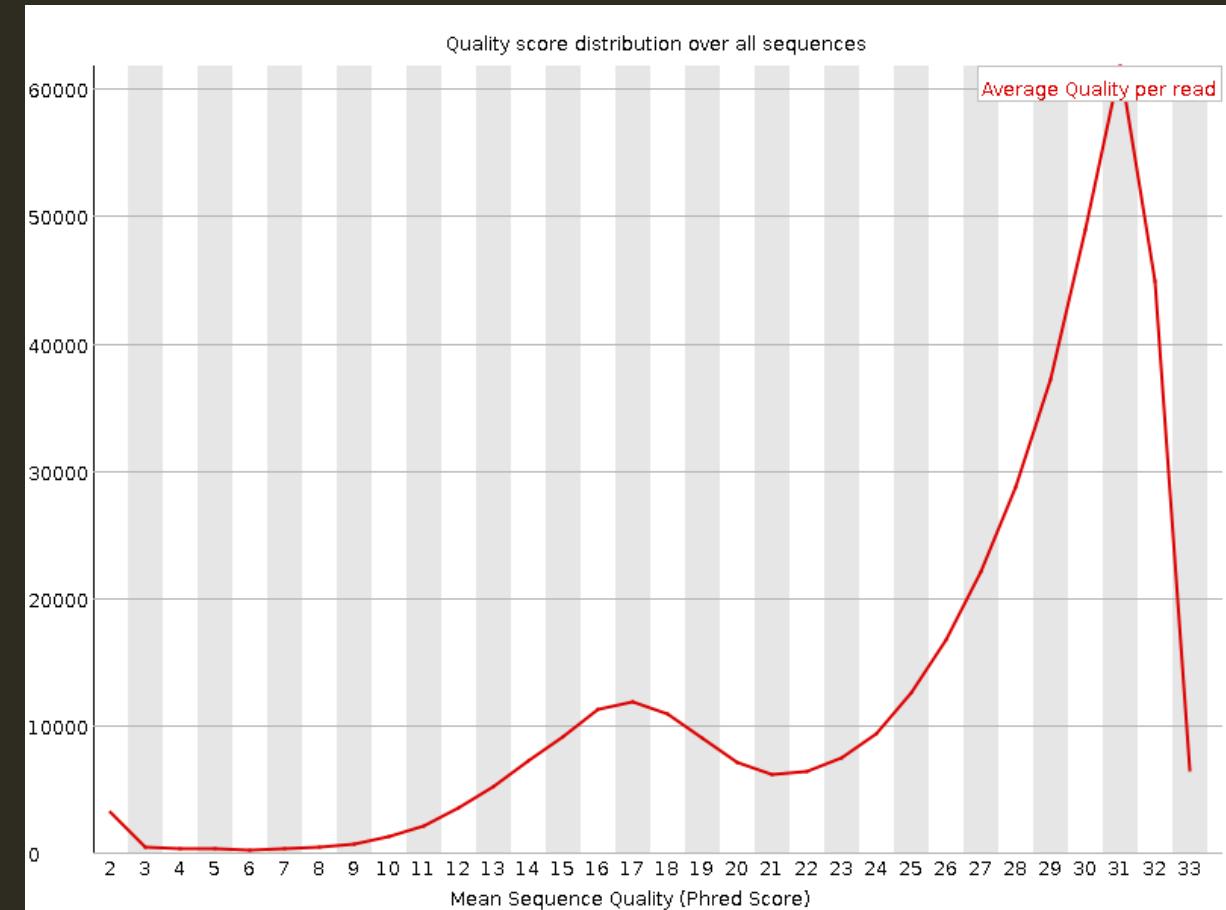
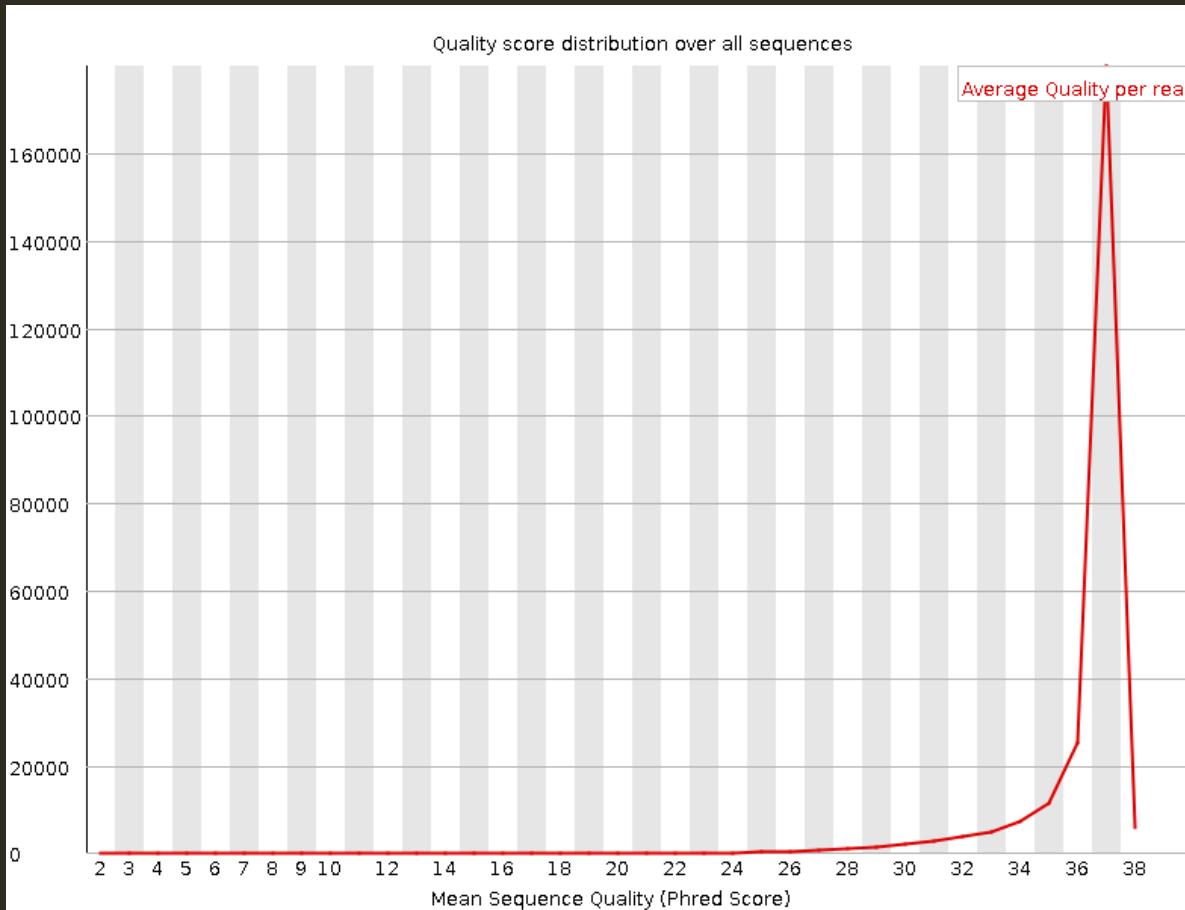
Basic Statistics

Measure	Value
Filename	bad_sequence.txt
File type	Conventional base calls
Encoding	Illumina 1.5
Total Sequences	395288
Sequences flagged as poor quality	0
Sequence length	40
%GC	47

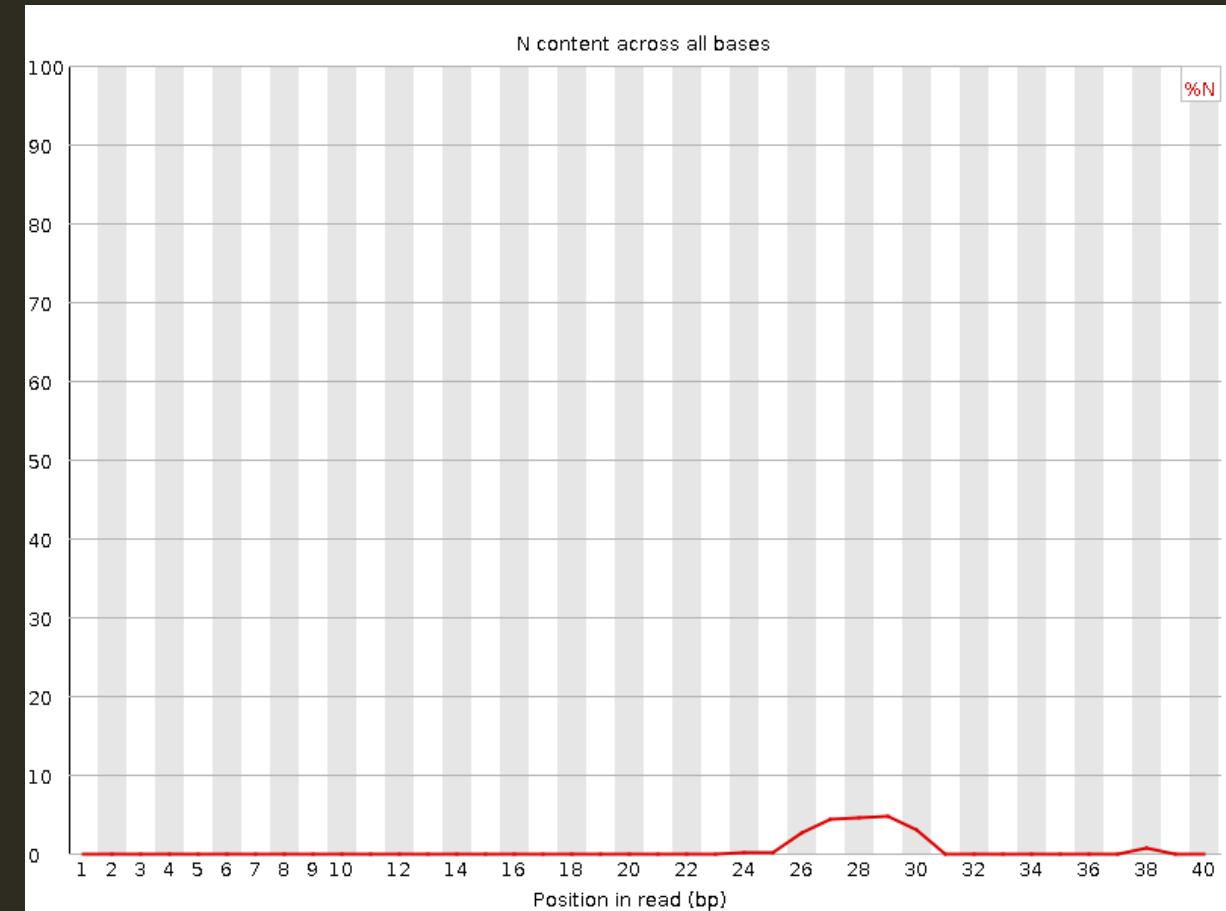
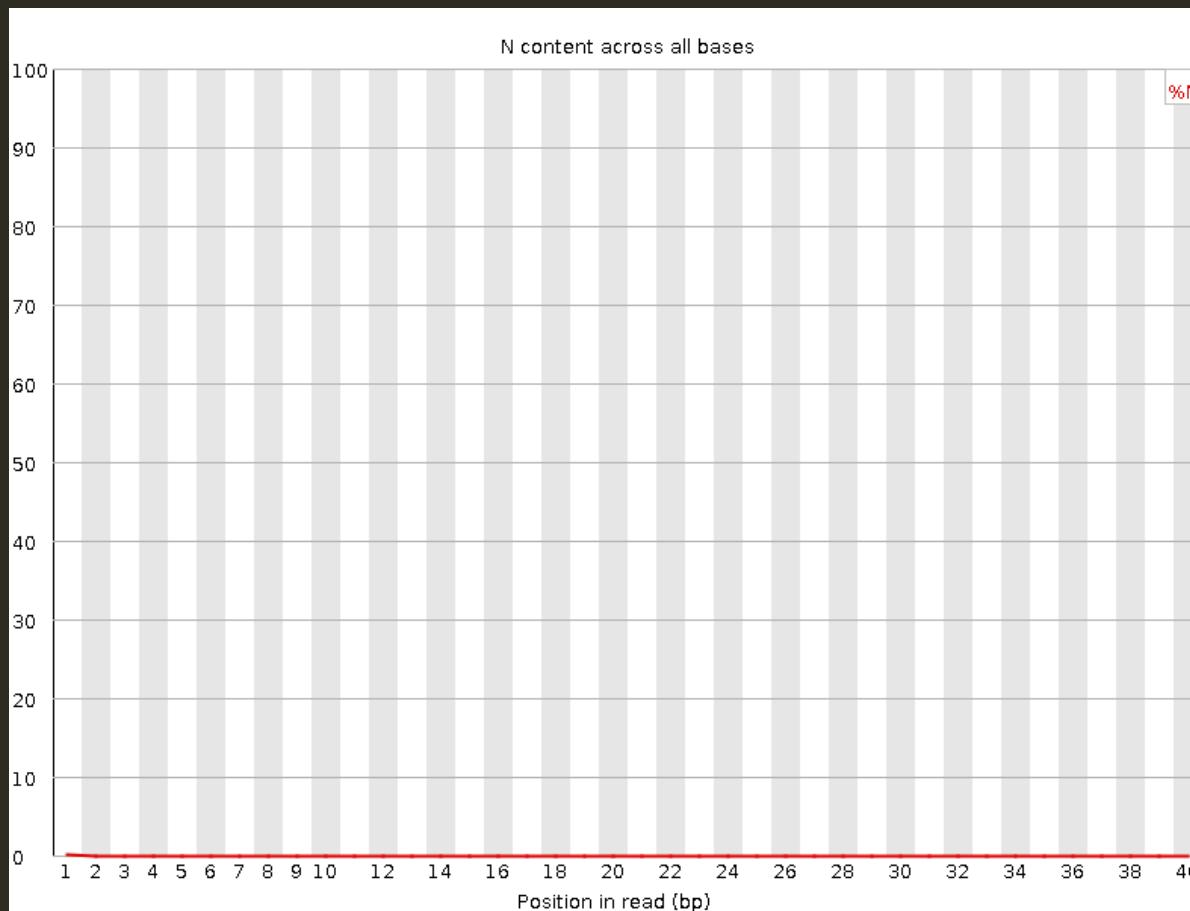
PER BASE SEQUENCE QUALITY



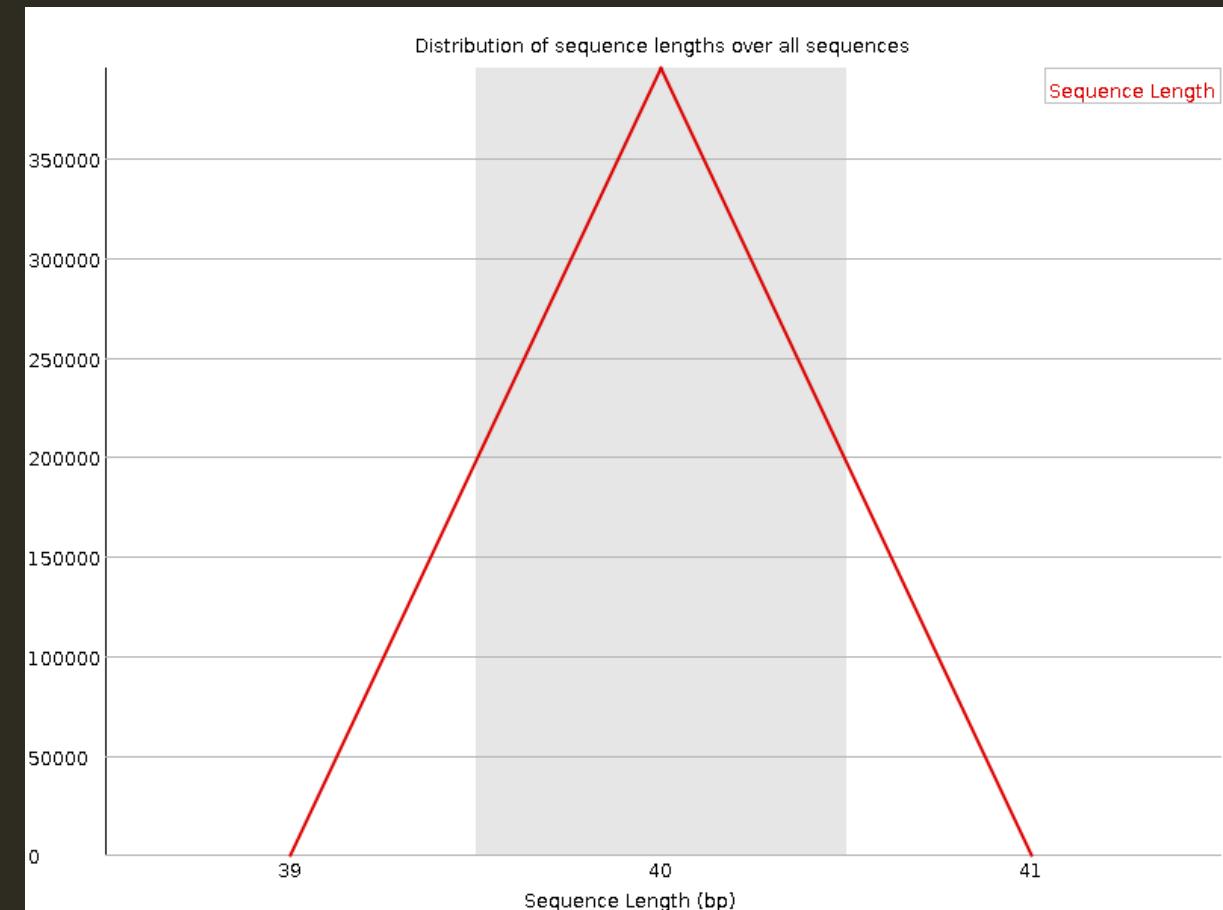
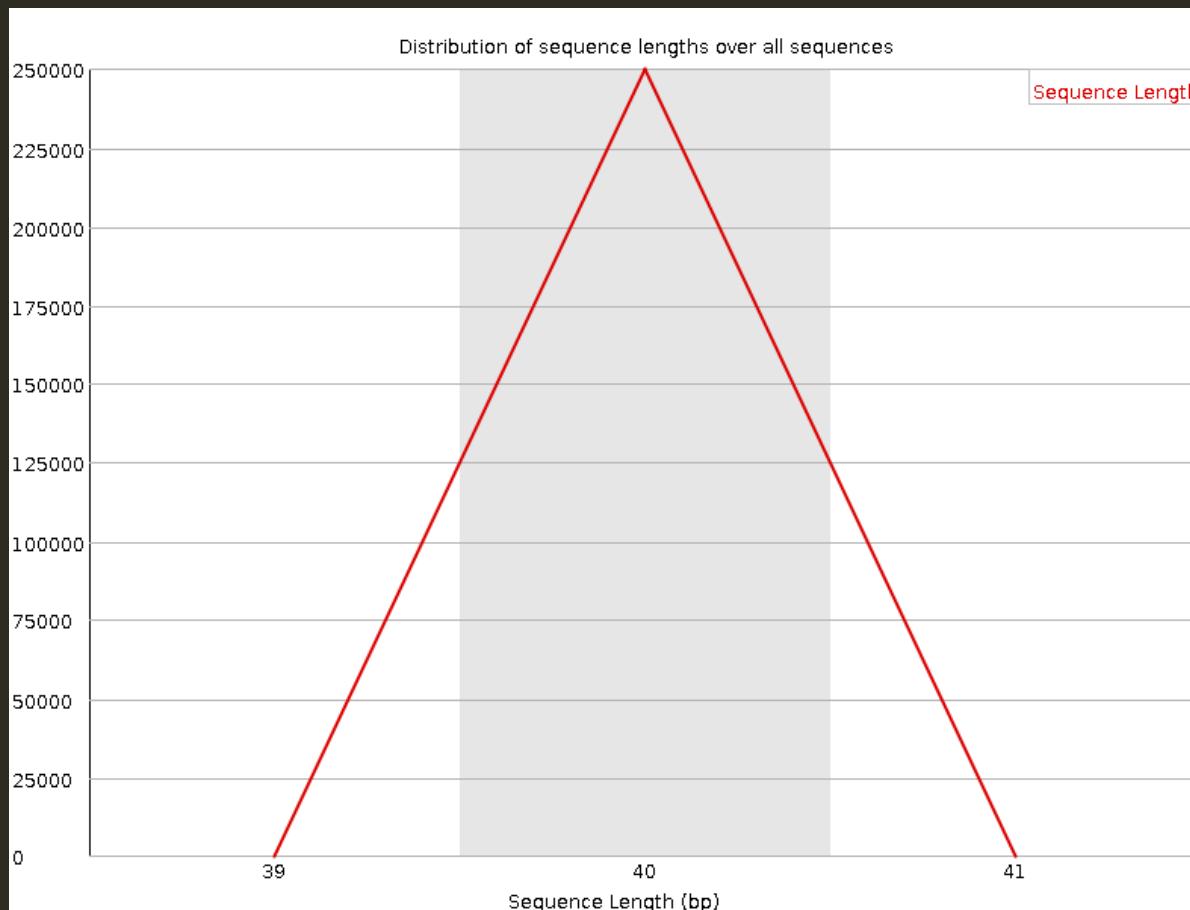
PER SEQUENCE QUALITY SCORES



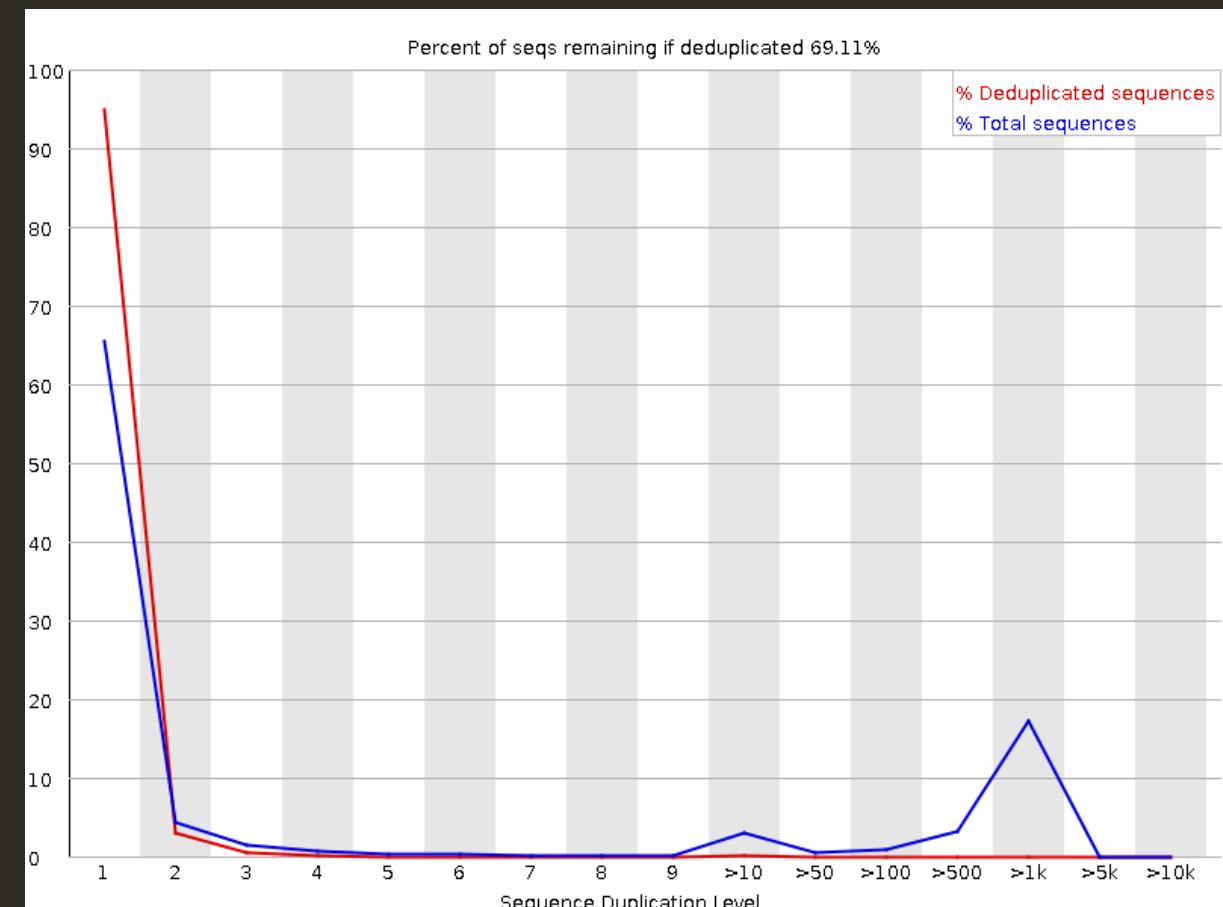
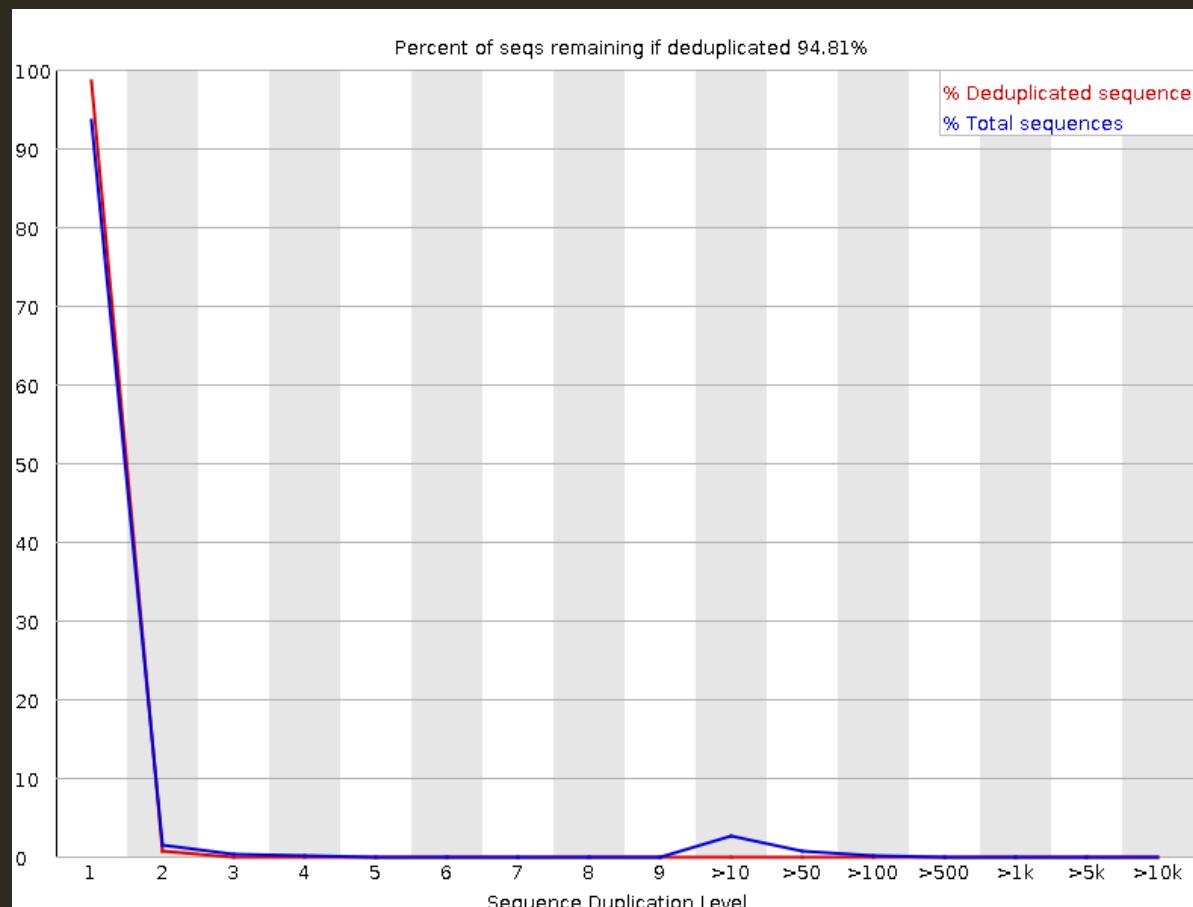
PER BASE N CONTENT



SEQUENCE LENGTH DISTRIBUTION



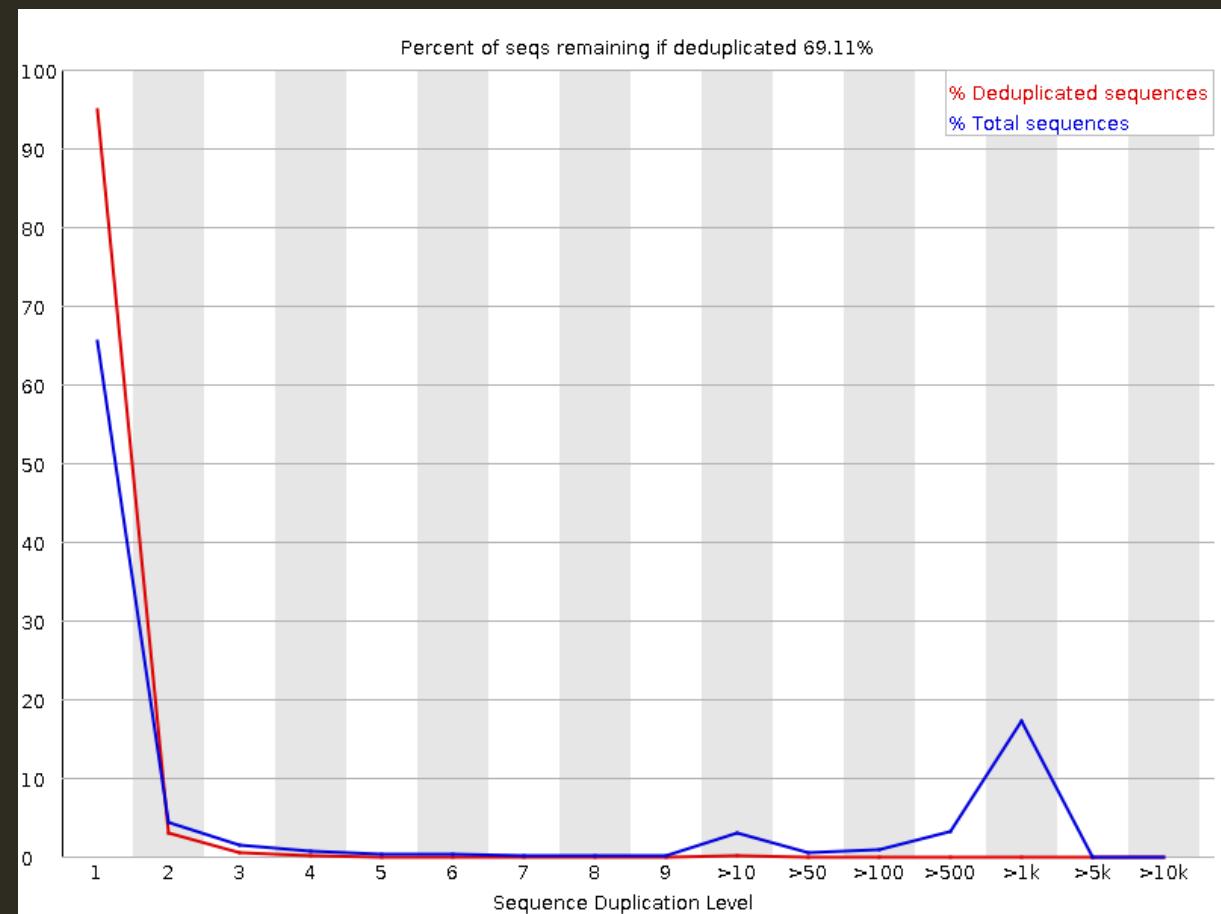
SEQUENCE DUPLICATION LEVELS



SEQUENCE DUPLICATION LEVELS

To cut down on the memory requirements for this module only sequences which first appear in the first **100,000 sequences** in each file are analysed, but this should be enough to get a good impression for the duplication levels in the whole file.

Because the duplication detection requires an exact sequence match over the whole length of the sequence, any reads over 75bp in length are truncated to **50bp** for the purposes of this analysis. Even so, longer reads are more likely to contain sequencing errors which will artificially increase the observed diversity and will tend to underrepresent highly duplicated sequences.



FASTQC - PRZYKŁADY

The main functions of FastQC are

- Import of data from BAM, SAM or FastQ files (any variant)
- Providing a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive appli

Documentation

A [copy of the FastQC documentation](#) is available for you to try before you buy (well download..).

Example Reports

- [Good Illumina Data](#)
- [Bad Illumina Data](#)
- [Adapter dimer contaminated run](#)
- [Small RNA with read-through adapter](#)
- [Reduced Representation BS-Seq](#)
- [PacBio](#)
- [454](#)



FASTQC - PRZYKŁADY

The main functions of FastQC are

- Import of data from BAM, SAM or FastQ files (any variant)
- Providing a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running

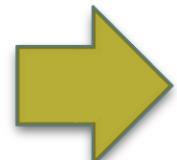
- [Basic Statistics](#)
- [Per base sequence quality](#)
- [Per sequence quality scores](#)
- [Per base sequence content](#)
- [Per sequence GC content](#)
- [Per base N content](#)
- [Sequence Length Distribution](#)
- [Sequence Duplication Levels](#)
- [Overrepresented sequences](#)
- [Adapter Content](#)

Documentation

A [copy of the FastQC documentation](#) is available for you to try before you buy.

Example Reports

- [Good Illumina Data](#)
- [Bad Illumina Data](#)
- [Adapter dimer contaminated run](#)
- [Small RNA with read-through adapter](#)
- [Reduced Representation BS-Seq](#)
- [PacBio](#)
- [454](#)



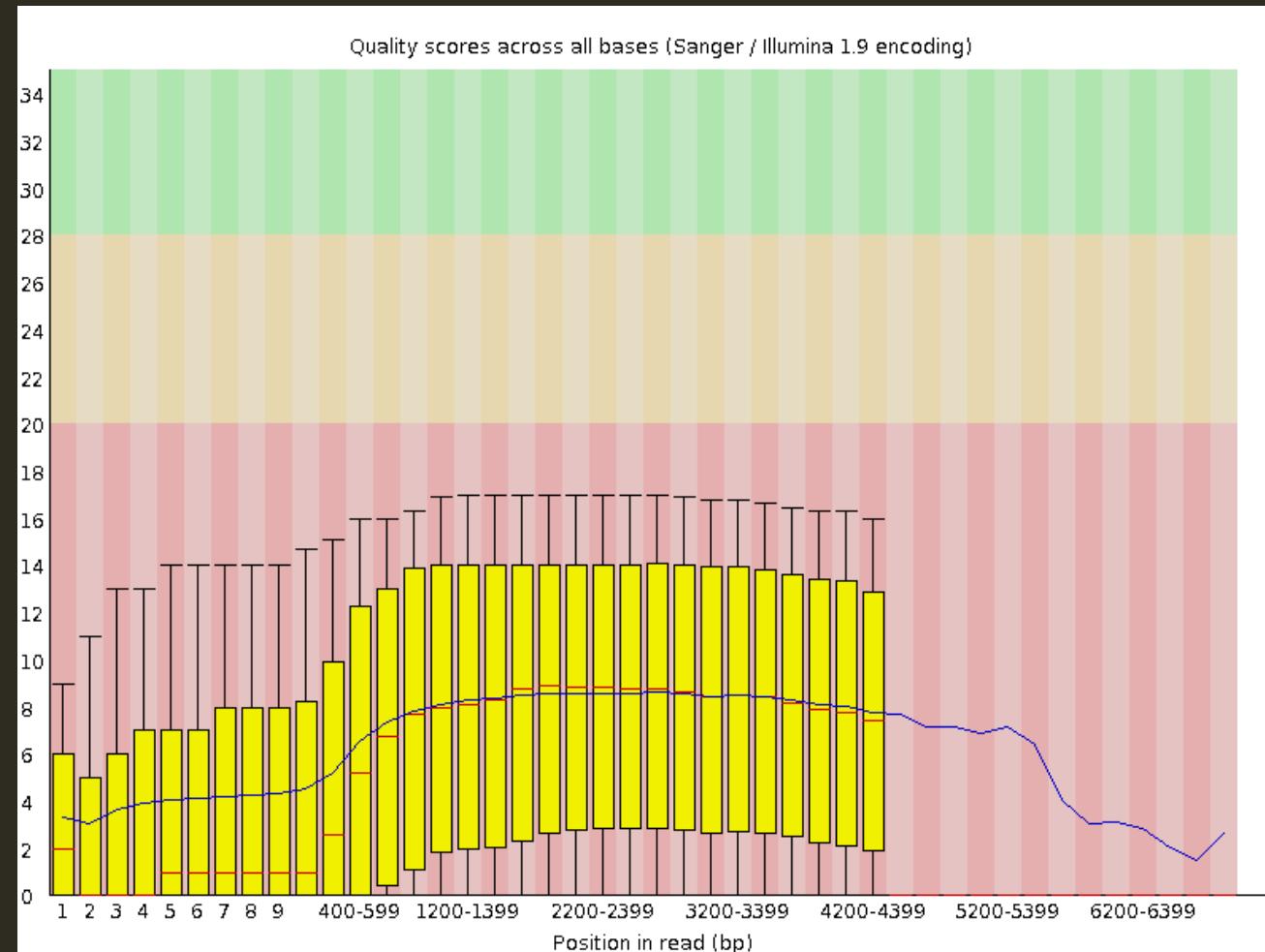
BASIC STATISTICS



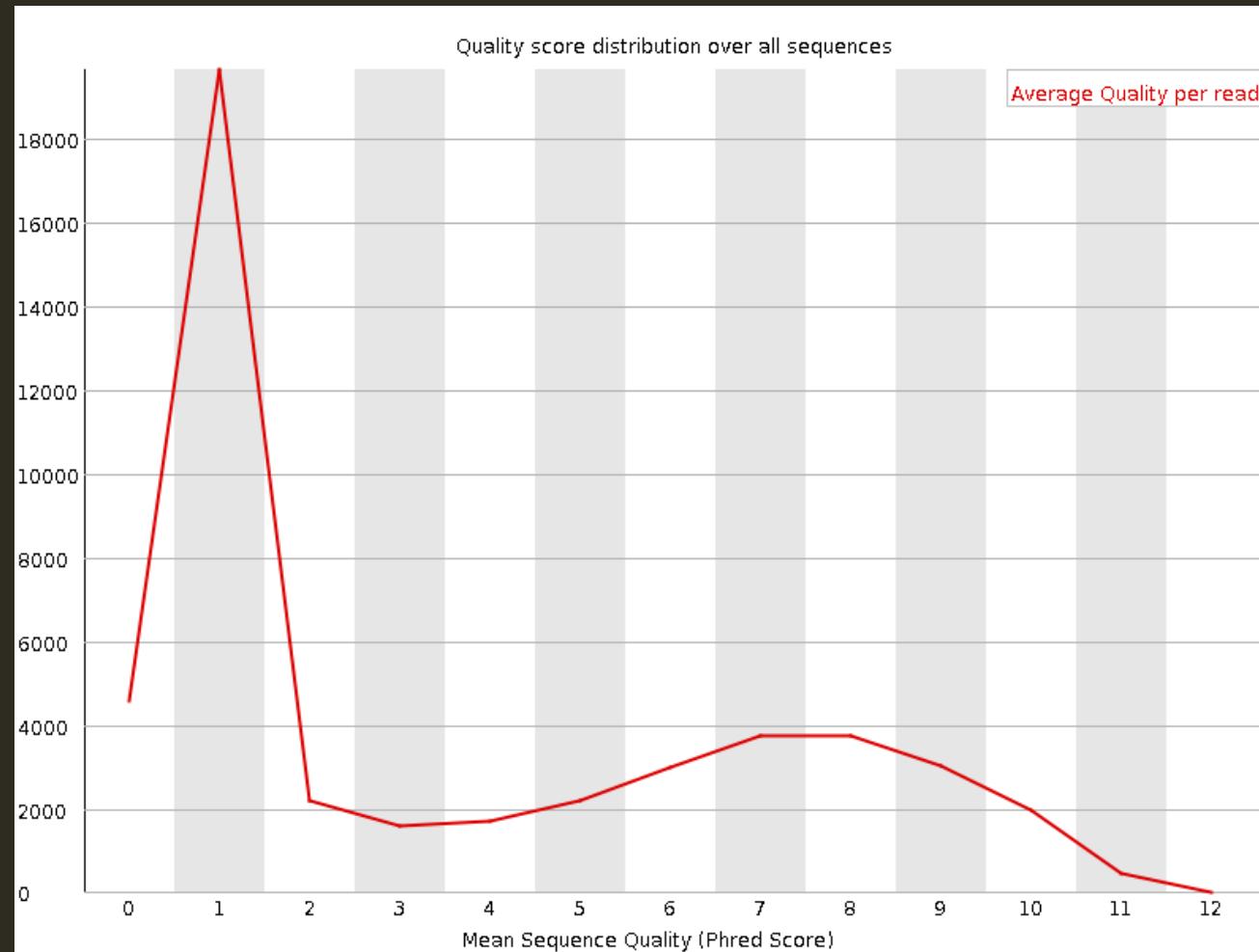
Basic Statistics

Measure	Value
Filename	pacbio_srr075104.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	48053
Sequences flagged as poor quality	0
Sequence length	82-6919
%GC	52

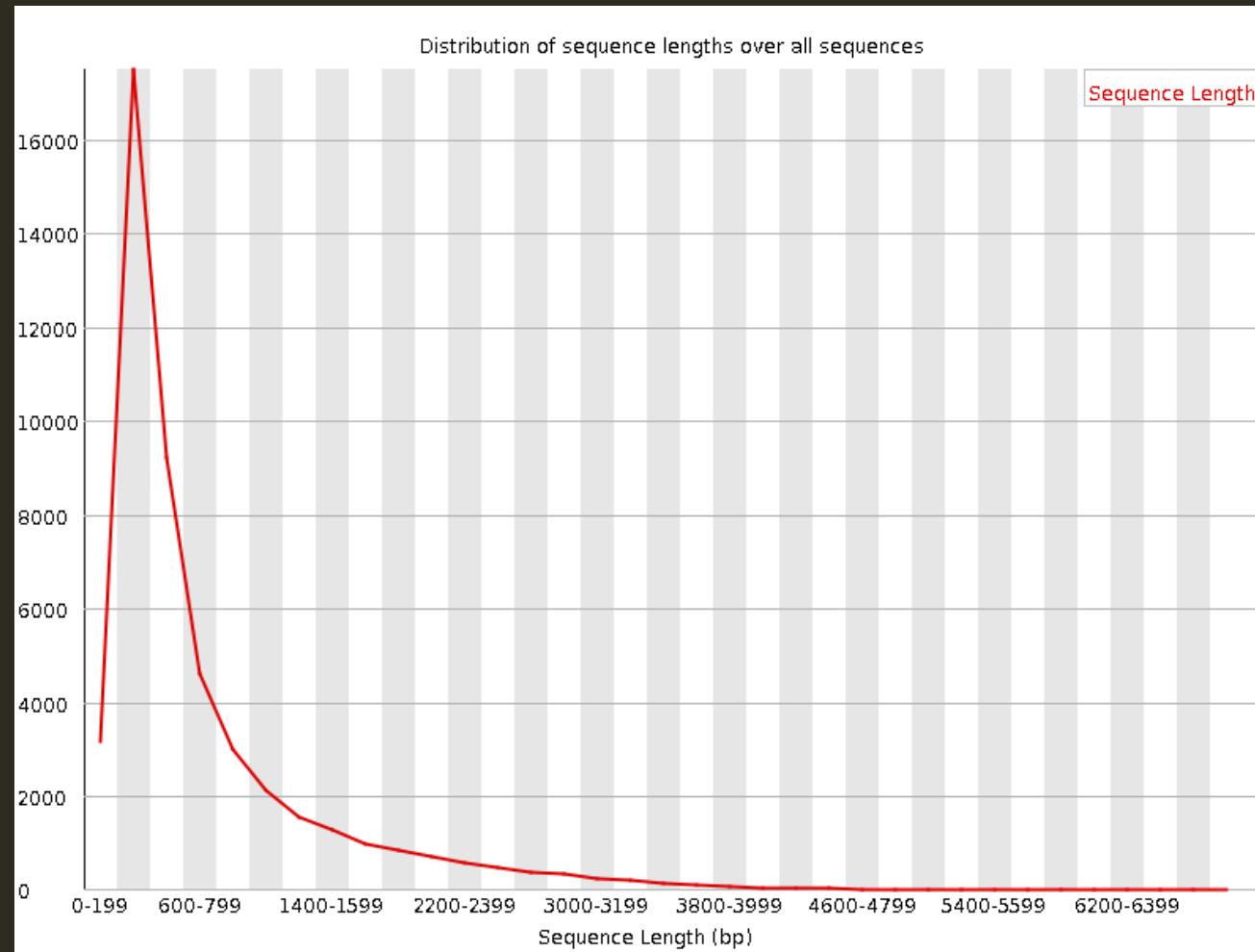
PER BASE SEQUENCE QUALITY



PER SEQUENCE QUALITY SCORES



SEQUENCE LENGTH DISTRIBUTION



EDYCJA SEKWENCJI

→ MOTYWACJA

Błędne dane mogą prowadzić do:

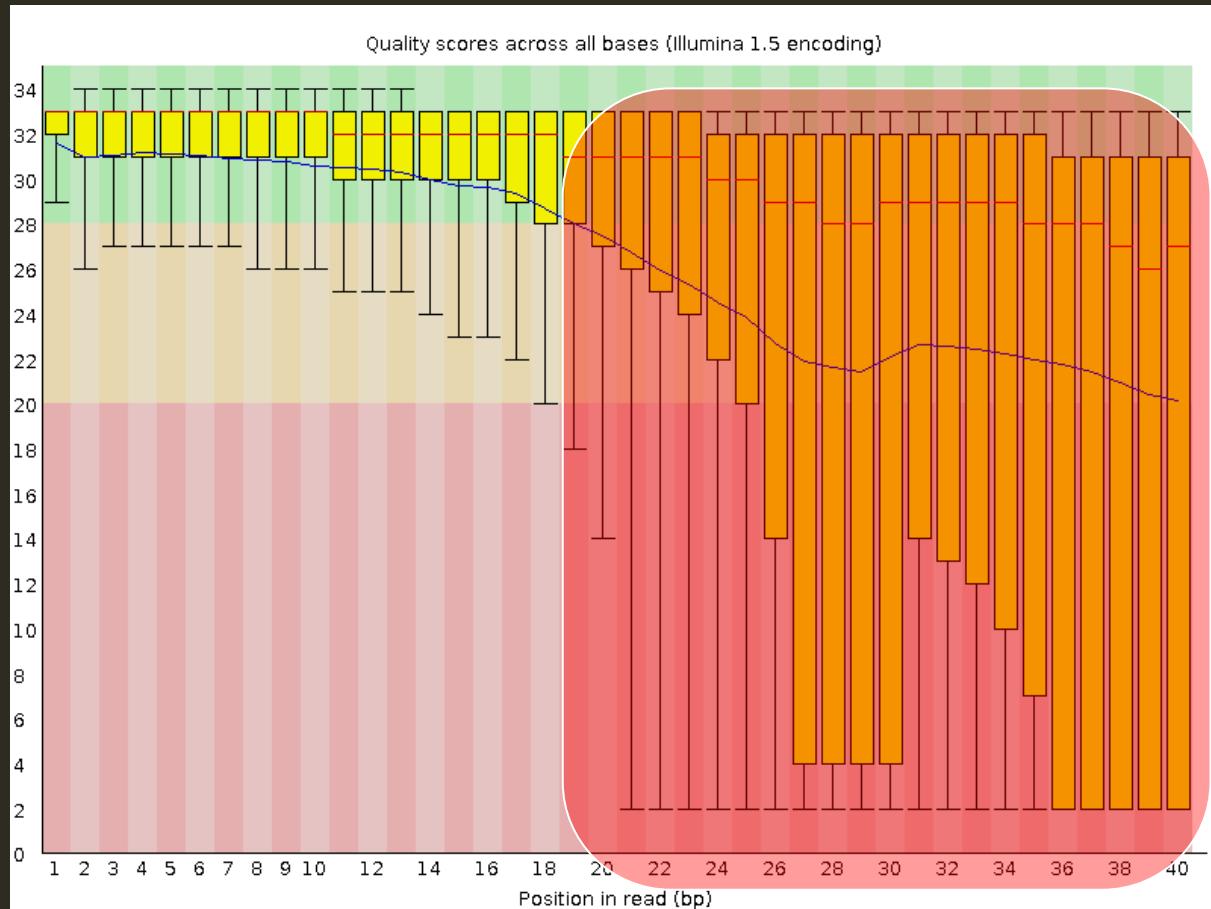
- wolniejszego działania oprogramowania
- generowania słabej jakości/niewłaściwych wyników

Czyszczenie danych:

- zwiększa średnią jakość sekwencji
- daje lepsze rezultaty przerównania
- redukuje rozmiar danych



```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
AGAAATGCCAGGCTAGATGAGTTACAATCNAGTATCAAGATAGGC  
+  
@@@FFDFFGHGHFDDGHHDDDD44#$%&,344+400/01234  
(...)
```



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

```
@HWI-1KL157:109:C448WACXX:7:1311:12007:37445 1:N:0:ACAGTG  
GTTAGCGCGCGGCTAGATGAGTTACAATCNAGTATCAAGATAGGAAAAAAA  
+  
@@@FFDFFGHHFDDDGHHHDDDD44#$%&,344+400/01234222211  
(...)
```

Oryginalny odczyt = 51 bp

1. Homopolimery?

TTAGCGCGCGGCTAGATGAGTTACAATCNAGTATCAAGATAGGAAAAAAA

2. Nieznane zasady?

TTAGCGCGCGGCTAGATGAGTTACAATCNAGTATCAAGATAGGAAAAAAA

3. Jakość poniżej 20?

TTAGCGCGCGGCTAGATGAGTTACAATCNAGTATCAAGATAGGAAAAAAA

Sekwencja po czyszczeniu = 26 bp

PRZYRÓWNANIE DO GENOMU REFERENCYJNEGO

ACTGGTGGGAA
ACTGGTGGGAA
GGTGGGAAAAAA
TGGGAAAAATT
GAAAAAAATTCA
GGGACTGATTCC
GACTGATTCCGA

AAAGGGAACCT
AAAGGGAACCT
GGGAACCTTTCT
GAACCTTCTTC
CCTTTCTTCGGA
AGAGAGATTTG
GAGAACCTTTCT

ACTGGTGGGAAAAATTCAAAAGGGAACCTTCTTGGAGCGGGACTGATTCCGAGAGAGA...

Genom referencyjny

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

PRZYRÓWNANIE DO GENOMU REFERENCYJNEGO



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

PRZYRÓWNANIE DO GENOMU REFERENCYJNEGO

Table 1 The characteristics of suffix array-based alignment to the reference genome software						
Name	Indexing	Output formats	PE mode	Gapped alignment	Supported platforms	Operating system
Bowtie	Genome	SAM	+	-	Illumina, ABI SOLiD	Linux, Macintosh, Windows, Solaris
Bowtie2	Genome	SAM	+	+	Illumina, 454, Ion Torrent	Linux, Macintosh, Windows
BWA	Genome	SAM	+	+	Illumina, 454, Ion Torrent	Linux

Surowe Dane

Kontrola jakości
Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

SAM (SEQUENCE ALIGNMENT/MAP FORMAT)

header section

```
@HD VN:1.0
@SQ SN:chr20 LN:62435964
@RG ID:L1 PU:SC_1_10 LB:SC_1 SM:NA12891
@RG ID:L2 PU:SC_2_12 LB:SC_2 SM:NA12891
read_28833_29006_6945 99 chr20 28833 20 10M1D25M = 28993 195 \
    AGCTTAGCTAGCTACCTATATCTTGGTCTTGGCCG <<<<<<<<<<<<<<<<|<9/,&,22,,<<< \
        NM:i:1 RG:Z:L1
read_28701_28881_323b 147 chr20 28834 30 35M      = 28701 -168 \
    ACCTATATCTTGGCCTTGGCCGATGCCGCCTTGCA <<<<, <<<<7, <<<<6, <<<<<<<<<7<<<< \
        MF:i:18 RG:Z:L2
```

alignment section

```

read_28833_29006_6945 99 chr20 28833 20 10M1D25M = 28993 195 \
    AGCTTAGCTAGCTACCTATATCTTGGTCTTGGCCG <<<<<<<<<<<<<<<,<9 / ,&,22,,<<< \
    NM:i:1 RG:Z:L1
read_28701_28881_323b 147 chr20 28834 30 35M      = 28701 -168 \
    ACCTATATCTTGGCCTTGGCGATGCGGCCTTGCA <<<<,<<<7,,<<<6,<<<<<<<<<7<<< \
    MF:i:18 RG:Z:L2

```

SAM/BAM

The alignment section consists of multiple TAB-delimited lines with each line describing an alignment. Each line is:

```

<QNAME> <FLAG> <RNAME> <POS> <MAPQ> <CIGAR> <MRNM> <MPOS> <ISIZE> <SEQ> <QUAL> \
[<TAG>:<VTYPE>:<VALUE> [ ... ]]

```

Field	Regular expression	Range	Description
QNAME	[^ \t\n\r]+		Query pair NAME if paired; or Query NAME if unpaired ²
FLAG	[0-9]+	[0,2 ¹⁶ -1]	bitwise FLAG (Section 2.2.2)
RNAME	[^ \t\n\r@=]+		Reference sequence NAME ³
POS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost POSition/coordinate of the clipped sequence
MAPQ	[0-9]+	[0,2 ⁸ -1]	MAPping Quality (phred-scaled posterior probability that the mapping position of this read is incorrect) ⁴
CIGAR	([0-9]+[MIDNSHP])+ *		extended CIGAR string
MRNM	[^ \t\n\r@]+		Mate Reference sequence NaMe; “=” if the same as <RNAME> ³
MPOS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost Mate POSition of the clipped sequence
ISIZE	-?[0-9]+	[-2 ²⁹ ,2 ²⁹]	inferred Insert SIZE ⁵
SEQ	[acgtnACGTN.=]+ *		query SEQuence; “=” for a match to the reference; n/N/ for ambiguity; cases are not maintained ^{6,7}
QUAL	[!-~]+ *	[0,93]	query QUALity; ASCII-33 gives the Phred base quality ^{6,7}
TAG	[A-Z][A-Z0-9]		TAG
VTYPE	[AifZH]		Value TYPE
VALUE	[^\t\n\r]+		match <VTYPE> (space allowed)

HWI-1KL157:58:D2FVAACXX:2:2313:3871:71331 147 Chr15 33794413 50 101M =
 33794252 -262
 GCTCAGCTTCTCACAGTCCAACTCTCACATCCATACATGACCACTGGAAAAACCATAGCCTGACTGGACGGACCT
 TTGTTAGAGGTTGCTAAAGACTG
 DBACCCDDCDDDC>DECDDC@;?3?3HAHGEJHHEFHD@CF=EDEEGDGG@IGF@HEHGCGIIJIG
 GDIHFGIJIJIEIIJIGG?FHHHFDDDD@@@ NM:i:2 AS:i:91 XS:i:83 RG:Z:D2FVAACXX_2

SAM/BAM

Field	Regular expression	Range	Description
QNAME	[^ \t\n\r]+		Query pair NAME if paired; or Query NAME if unpaired ²
FLAG	[0-9]+	[0,2 ¹⁶ -1]	bitwise FLAG (Section 2.2.2)
RNAME	[^ \t\n\r@=]+		Reference sequence NAME ³
POS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost POSition/coordinate of the clipped sequence
MAPQ	[0-9]+	[0,2 ⁸ -1]	MAPping Quality (phred-scaled posterior probability that the mapping position of this read is incorrect) ⁴
CIGAR	([0-9]+[MIDNSHP])+ *		extended CIGAR string
MRNM	[^ \t\n\r@]+		Mate Reference sequence NaMe; “=” if the same as <RNAME> ³
MPOS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost Mate POSition of the clipped sequence
ISIZE	-?[0-9]+	[-2 ²⁹ ,2 ²⁹]	inferred Insert SIZE ⁵
SEQ	[acgttnACGTN.=]+ *		query SEQuence; “=” for a match to the reference; n/N/ for ambiguity; cases are not maintained ^{6,7}
QUAL	[!-~]+ *	[0,93]	query QUALity; ASCII-33 gives the Phred base quality ^{6,7}
TAG	[A-Z][A-Z0-9]		TAG
VTYPE	[AifZH]		Value TYPE
VALUE	[^\t\n\r]+		match <VTYPE> (space allowed)

HWI-1KL157:58:D2FVAACXX:2:2313:3871:71331 147 Chr15 33794413 50 101M =
33794252 -262
GCTCAGCTTCTCACAGTCCAACTCTCACATCCATACATGACCACTGGAAAAACCATAGCCTGACTGGACGGACCT
TTGTTAGAGGTTGCTAAAGACTG
DBACCCDDCDDDC>DECDDC@;?3?3HAHGEJHHEFHD@CF=EDEEGDGG@IGF@HEHGCGIIJIG
GDIHFGIJIJIEIIJIGG?FHHHFDDDD@@@ NM:i:2 AS:i:91 XS:i:83 RG:Z:D2FVAACXX_2

SAM/BAM

<https://broadinstitute.github.io/picard/explain-flags.html>

Decoding SAM flags

This utility makes it easy to identify what are the properties of a read based on its SAM flag value, or conversely, to find what the SAM Flag value would be for a given combination of properties.

To decode a given SAM flag value, just enter the number in the field below. The encoded properties will be listed under Summary below, to the right.

SAM Flag: [Explain](#)

[Switch to mate](#) Toggle first in pair / second in pair

Find SAM flag by property:

To find out what the SAM flag value would be for a given combination of properties, tick the boxes for those that you'd like to include. The flag value will be shown in the SAM Flag field above.

- read paired
- read mapped in proper pair
- read reverse strand

Summary:

- read paired (0x1)
- read mapped in proper pair (0x2)
- read reverse strand (0x10)
- second in pair (0x80)

HWI-1KL157:58:D2FVAACXX:2:2313:3871:71331 147 Chr15 33794413 50 101M =
 33794252 -262
 GCTCAGCTTCTCACAGTCCAACTCTCACATCCATACATGACCACTGGAAAAACCATAGCCTGACTGGACGGACCT
 TTGTTAGAGGTTGCTAAAGACTG
 DBACCCDDCDDDC>DECDDC@;?3?3HAHGEJHHEFHD@CF=EDEEGDGG@IGF@HEHGCGIIJIG
 GDIHFGIJIJIEIIJIGG?FHHHFDDDD@@@ NM:i:2 AS:i:91 XS:i:83 RG:Z:D2FVAACXX_2

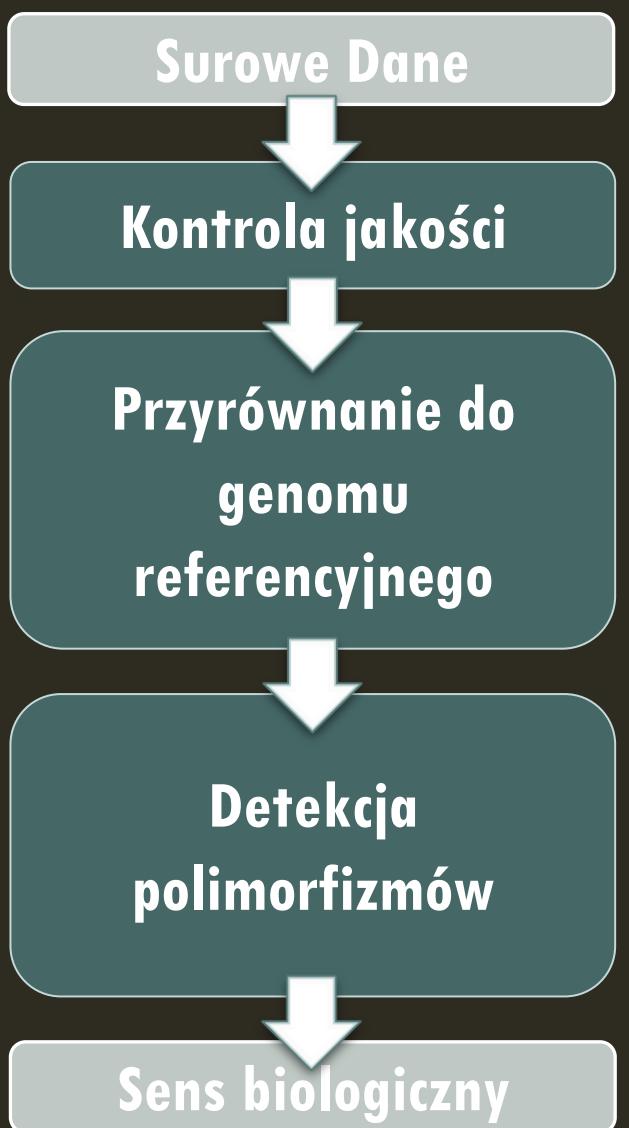
SAM/BAM

Field	Regular expression	Range	Description
QNAME	[^ \t\n\r]+		Query pair NAME if paired; or Query NAME if unpaired ²
FLAG	[0-9]+	[0,2 ¹⁶ -1]	bitwise FLAG (Section 2.2.2)
RNAME	[^ \t\n\r@=]+		Reference sequence NAME ³
POS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost POSition/coordinate of the clipped sequence
MAPQ	[0-9]+	[0,2 ⁸ -1]	MAPping Quality (phred-scaled posterior probability that the mapping position of this read is incorrect) ⁴
CIGAR	([0-9]+[MIDNSHP])+ *		extended CIGAR string
MRNM	[^ \t\n\r@]+		Mate Reference sequence NaMe; “=” if the same as <RNAME> ³
MPOS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost Mate POSition of the clipped sequence
ISIZE	-?[0-9]+	[-2 ²⁹ ,2 ²⁹]	inferred Insert SIZE ⁵
SEQ	[acgttnACGTN.=]+ *		query SEQuence; “=” for a match to the reference; n/N/ for ambiguity; cases are not maintained ^{6,7}
QUAL	[!-~]+ *	[0,93]	query QUALity; ASCII-33 gives the Phred base quality ^{6,7}
TAG	[A-Z][A-Z0-9]		TAG
VTYPE	[AifZH]		Value TYPE
VALUE	[^\t\n\r]+		match <VTYPE> (space allowed)

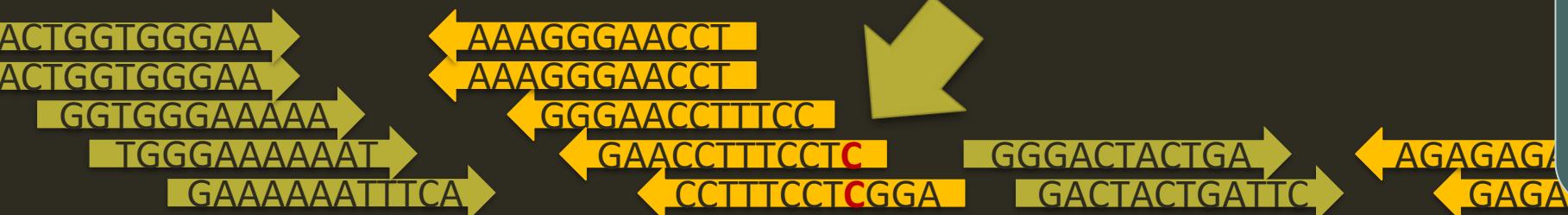
DETEKCJA POLIMORFIZMÓW

- Single Nucleotide Polymorphisms
- Insertions/Deletions
- Copy Number Variations
- Loss of Heterozygosity
- Inversions
- Translocations

SNP
INDEL
CNV
LOH
INV
TRANS



DETEKCJA POLIMORFIZMÓW



Genom referencyjny

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

DETEKCJA POLIMORFIZMÓW → SNP

Name	Input formats	Output formats	Realignment	Recalibration	Single-sample	Multi-sample	Called variants	Operating system
Atlas-SNP2	BAM	VCF	–	–	+	–	SNPs	Linux
GATK (UnifiedGenotyper)	BAM	VCF	+	+	+	+	SNPs, InDels	Linux
SAMtools (samtools mpileup)	BAM	VCF	+	+	+	+	SNPs, InDels	Linux
SNVer	BAM	VCF	–	–	+	–	SNPs, InDels	Linux, Macintosh, Windows
SOAPsnp	SOAP	Text format, GLFv2	–	+	+	–	SNPs	Linux
VarScan2	Pileup/mpileup	Text format, VCF	–	–	+	+	SNPs, InDels, CNA	Linux, Macintosh, Windows

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

```
##fileformat=VCFv4.2
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1>Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1>Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A>Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1>Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0>Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0>Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1>Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1>Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1>Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2>Type=Integer,Description="Haplotype Quality">
```

CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00003
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	0/0:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2
20	1234567	microsat1	GTC	G,GTCT	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3

polimorphisms (body)

VCF

#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE

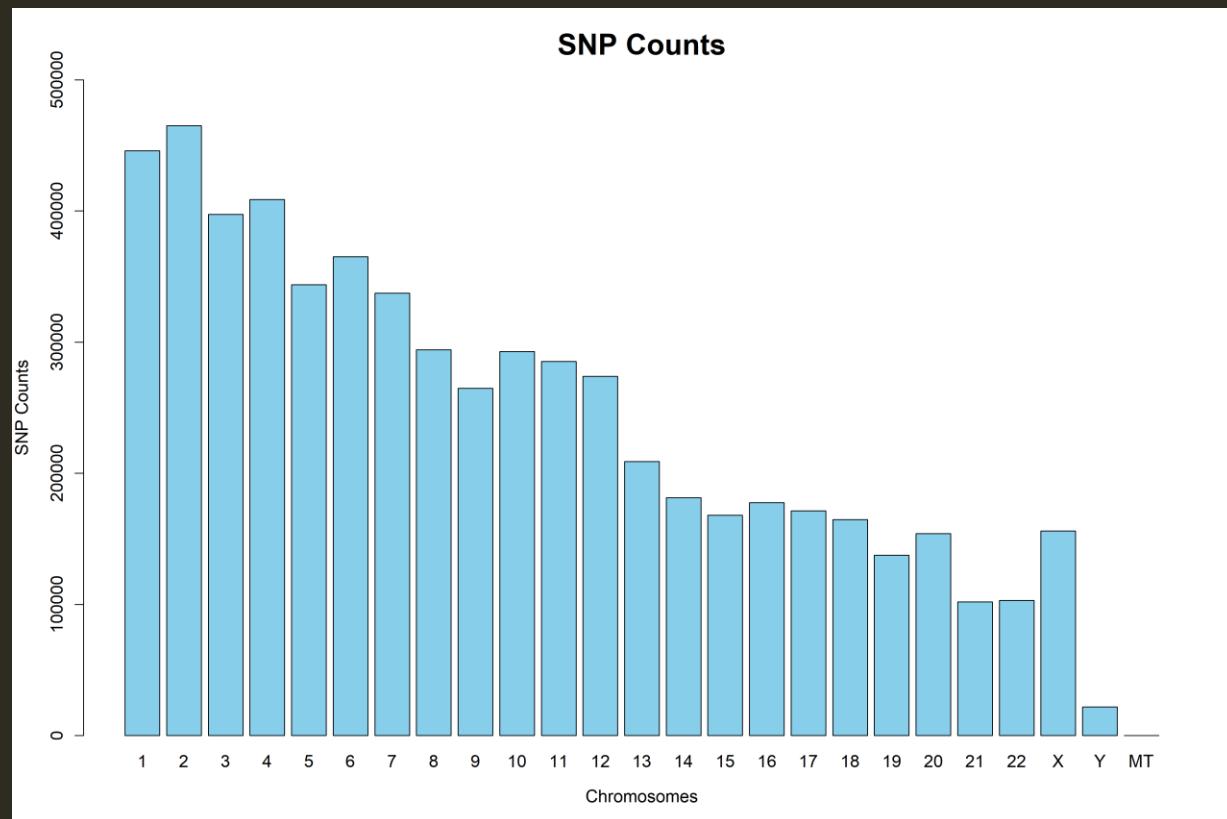
Chr1 238 . C T 48 . DP=8;VDB=4.789490e-02;RPB=-1.551181e+00;AF1=0.5;AC1=1;DP4=4,1,3,0;MQ=44;FQ=51; PV4=1,1,0.16,0.41
GT:PL:GQ 0/1:78,0,135:81

```
##fileformat=VCFv4.2
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1>Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1>Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A>Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1>Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0>Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0>Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1>Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1>Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1>Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2>Type=Integer,Description="Haplotype Quality">
```

#CHROM POS ID REF ALT QUAL FILTER INFO							FORMAT	NA00001	NA00002	NA00003	
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	0/0:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2
20	1234567	microsat1	GTC	G,GTCT	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3

DETEKCJA POLIMORFIZMÓW → SNP

Homo sapiens



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

SENS BIOLOGICZNY

→ FUNKCJONALNA ADNOTACJA

Variant Effect Predictor

VEP determines the effect of your variants (SNPs, insertions, deletions, CNVs or structural variants) on genes, transcripts, and protein sequence, as well as regulatory regions. Simply input the coordinates of your variants and the nucleotide changes to find out the:

- genes and transcripts affected by the variants
- location of the variants (e.g. upstream of a transcript, in coding sequence, in non-coding RNA, in regulatory regions)
- consequence of your variants on the protein sequence (e.g. stop gained, missense, stop lost, frameshift)
- known variants that match yours, and associated minor allele frequencies from the 1000 Genomes Project
- SIFT and PolyPhen scores for changes to protein sequence
- ... And [more!](#)



→ www.ensembl.org/Tools/VEP

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

SENS BIOLOGICZNY

→ FUNKCJONALNA ADNOTACJA

The screenshot shows the Ensembl VEP tool interface. The top navigation bar includes links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. The main menu has 'Species' and 'VEP' dropdowns. A sidebar on the left lists 'Web Tools' such as Web Tools, BLAST/BLAT, Variant Effect Predictor (which is selected), File Chameleon, Assembly Converter, and ID History Converter. Other buttons include 'Configure this page', 'Custom tracks', 'Export data', 'Share this page', and 'Bookmark this page'. The main content area is titled 'Variant Effect Predictor' and displays information for 'Human GRCh37'. It says, 'If you are looking for VEP for Human GRCh37, please go to [GRCh37 website](#)'. Below this, there are fields for 'Species' (set to Human (Homo sapiens)) and 'Assembly' (set to GRCh38.p7). There is also a field for 'Name for this job (optional)' and a text area for 'Either paste data:' containing a list of variant IDs and descriptions. At the bottom, examples are given for Ensembl default, VCF, Variant identifiers, and HGVS notations.

Variant Effect Predictor ?

i VEP for Human GRCh37

If you are looking for VEP for Human GRCh37, please go to [GRCh37 website](#).

Species: Human (Homo sapiens)

Assembly: GRCh38.p7

Name for this job (optional):

Either paste data:

```
1 182712 182712 A/C 1
2 265023 265023 C/T 1
3 319781 319781 A/- 1
19 110748 110747 -/T 1
1 160283 471362 DUP 1
1 1385015 1387562 DEL 1
```

Examples: [Ensembl default](#), [VCF](#), [Variant identifiers](#), [HGVS notations](#).
NB: pileup format no longer supported

Surowe Dane

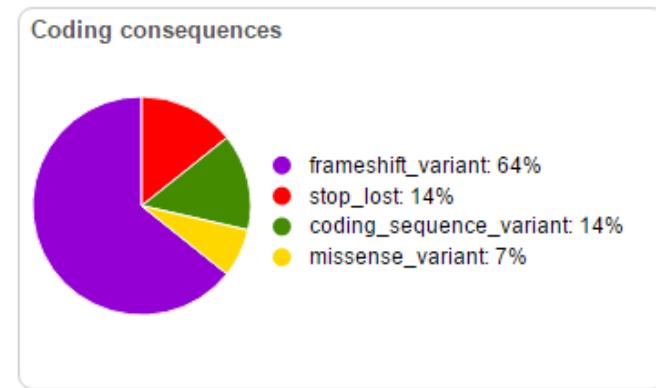
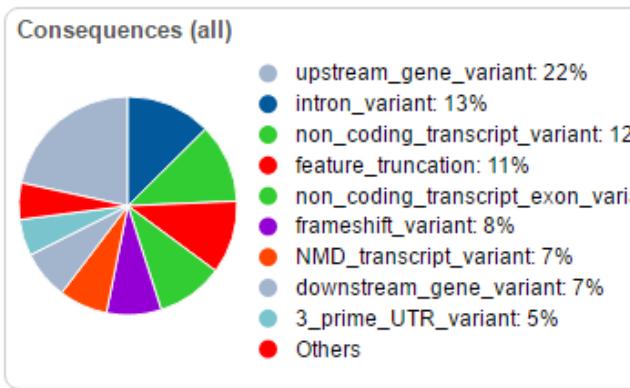
Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

Category	Count
Variants processed	6
Variants filtered out	0
Novel / existing variants	5 (83.3) / 1 (16.7)
Overlapped genes	10
Overlapped transcripts	70
Overlapped regulatory features	1



Results preview

Navigation
Filters
Download

Page: 1 of 1 | Show: All variants

Uploaded variant is defined

All: VCF VEP TXT
BioMart: Variants Genes

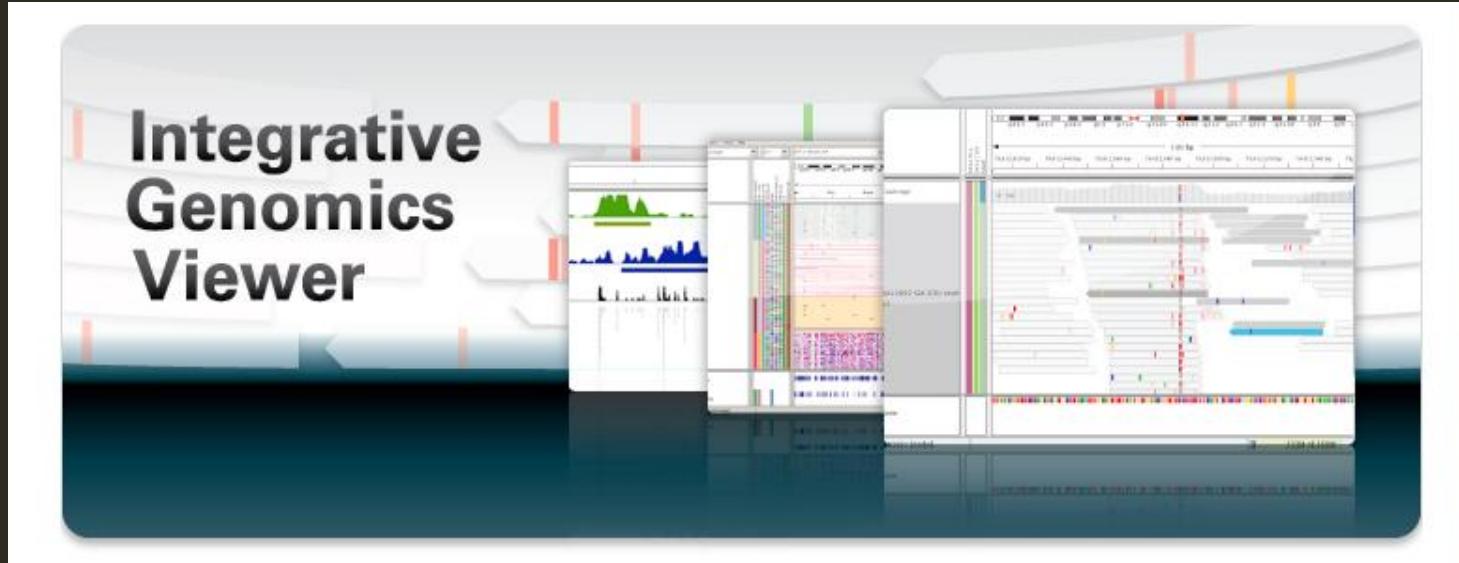
Show/hide columns

Uploaded variant	Location	Allele	Consequence	Impact	Symbol	Gene	Feature type	Feature	Biotype
1_160283_duplication	1:160282-160282	duplication	upstream_gene_variant	MODIFIER	RNU6-1100P	ENSG0000222623	Transcript	ENST0000410691	snRNA
1_160283_duplication	1:160282-160282	duplication	non_coding_transcript_exon_variant, intron_variant, non_coding_transcript_variant	MODIFIER	RP11-34P13.13	ENSG0000241860	Transcript	ENST0000466557	lincRNA
1_160283_duplication	1:160282-160282	duplication	non_coding_transcript_exon_variant, intron_variant, non_coding_transcript_variant	MODIFIER	RP11-34P13.13	ENSG0000241860	Transcript	ENST0000491962	lincRNA
1_160283_duplication	1:160282-160282	duplication	non_coding_transcript_exon_variant, intron_variant, non_coding_transcript_variant	MODIFIER	RP11-34P13.9	ENSG0000241599	Transcript	ENST0000496488	lincRNA
1_182712_A/C	1:182712-182712	C	downstream_gene_variant	MODIFIER	F0538757.1	ENSG0000279457	Transcript	ENST0000623083	protein_coding
1_182712_A/C	1:182712-182712	C	downstream_gene_variant	MODIFIER	F0538757.1	ENSG0000279457	Transcript	ENST0000623834	protein_coding
1_182712_A/C	1:182712-182712	C	missense_variant	MODERATE	F0538757.2	ENSG0000279928	Transcript	ENST0000624431	protein_coding
1_182712_A/C	1:182712-182712	C	downstream_gene_variant	MODIFIER	F0538757.1	ENSG0000279457	Transcript	ENST0000624735	protein_coding

WIZUALIZACJA

Wizualizacja w programie IGV

<http://software.broadinstitute.org/software/igv/>



Surowe Dane

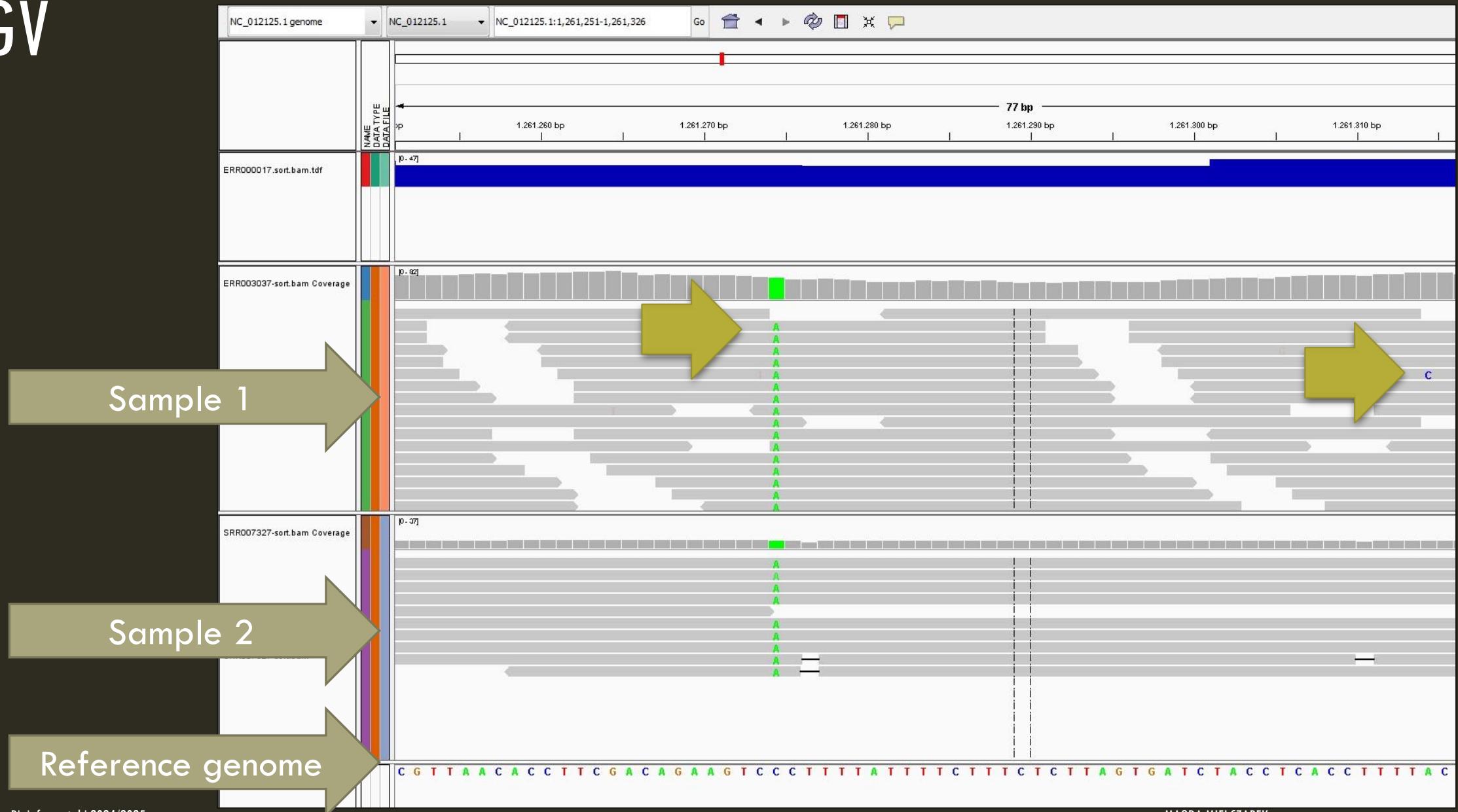
Kontrola jakości

Przyrównanie do
genomu
referencyjnego

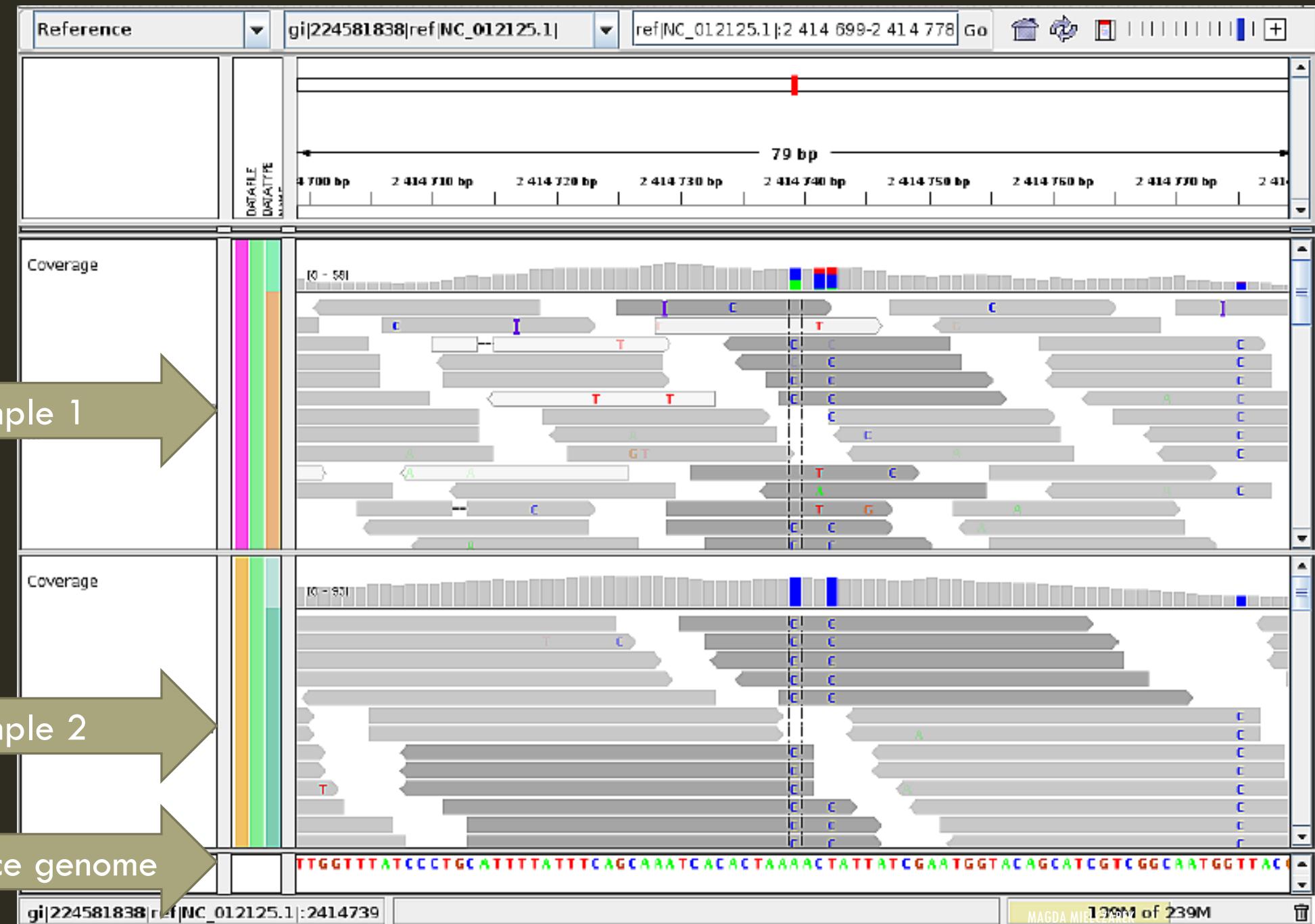
Detekcja
polimorfizmów

Sens biologiczny

IGV



IGV



SENS BIOLOGICZNY → MEDYCZYNA PERSONALIZOWANA

doi: [10.2217/pme.14.34](https://doi.org/10.2217/pme.14.34)

The road from next-generation sequencing to personalized medicine

Manuel L. Gonzalez-Garay

[Author information](#) ► [Copyright and License information](#) ►

See other articles in PMC that cite the published article.

Abstract

Go to: [View Article Online](#)

Moving from a traditional medical model of treating pathologies to an individualized predictive and preventive model of personalized medicine promises to reduce the healthcare cost on an overburdened and overwhelmed system. Next-generation sequencing (NGS) has the potential to accelerate the early detection of disorders and the identification of pharmacogenetics markers to customize treatments. This review explains the historical facts that led to the development of NGS along with the strengths and weakness of NGS, with a special emphasis on the analytical aspects used to process NGS data. There are solutions to all the steps necessary for performing NGS in the clinical context where the majority of them are very efficient, but there are some crucial steps in the process that need immediate attention.

Keywords: CADD, functional prediction program, genomics, GWAVA, NGS, personalized medicine, workflow management system

The current medical model focuses on the detection and treatment of pathologies. Treating disorders, especially on advanced states, is very expensive for patients and society in general. Screening for five of the most common disorders in the USA (cardiovascular disorders, stroke, cancer, chronic obstructive pulmonary disease and diabetes) could protect millions of lives and reduce the healthcare deficit [1]. Tailoring drug therapies by practicing personalized medicine (PM) has the potential to improve treatment of cancer and save lives by preventing drug-related adverse effects. A new technology, next-generation sequencing (NGS), has the potential to accelerate the early detection of disorders and to detect pharmacogenetics markers to customize treatments [2].

Use of NGS to diagnose human disorders

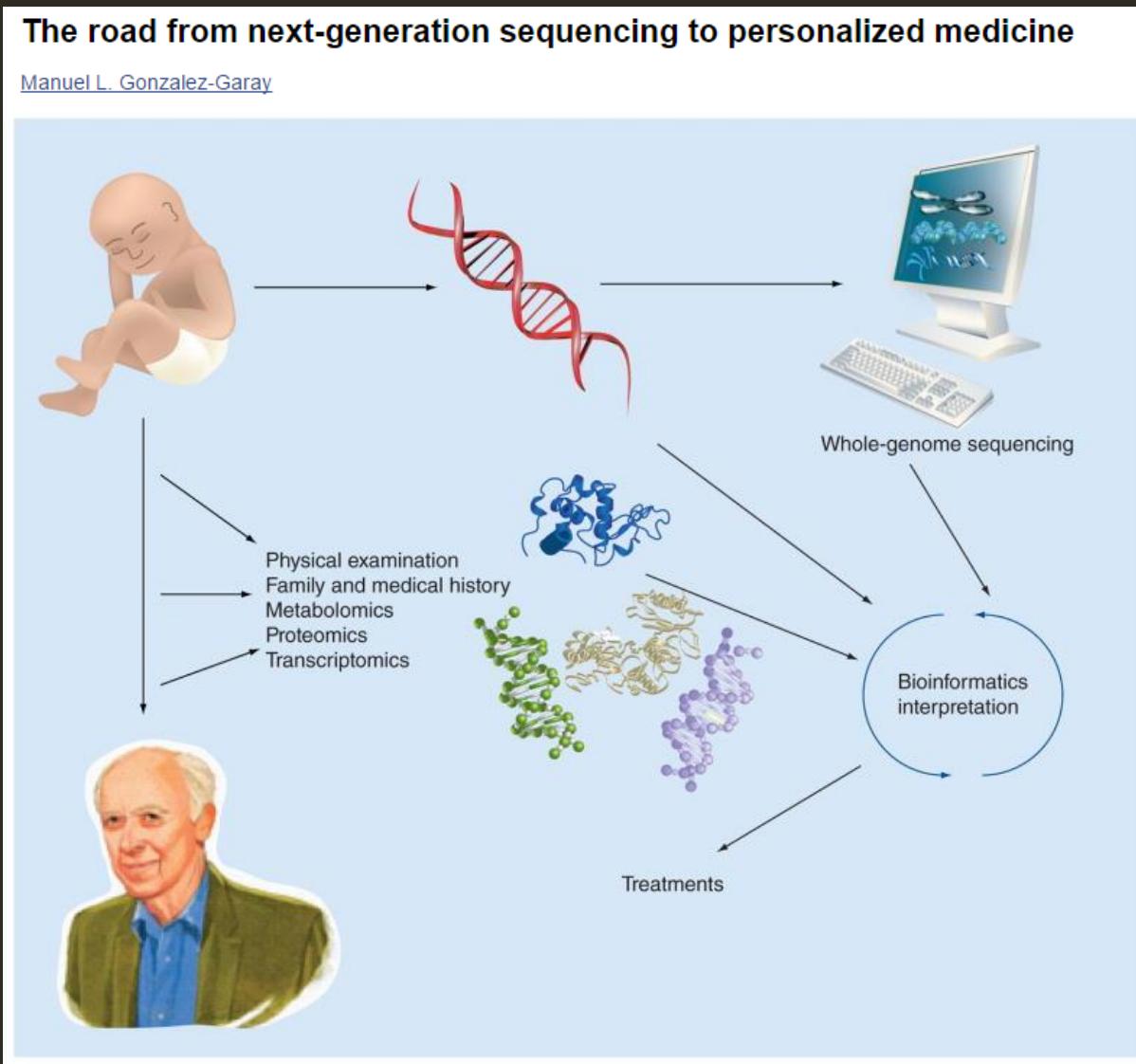
Go to: [View Article Online](#)

One of the major concerns of medical diagnosis is to identify genes and mutations responsible for human diseases. Early identification of causative mutations enables the early detection of a myriad of disorders. We are living in an age of high healthcare cost. Early detection of genetic disorders, carrier status, genetic predispositions for cancer and cardiovascular disease could potentially reduce the healthcare cost.

The first proof of concept that the NGS technology could be used to detect genetic disorders was provided by Shendure's group on September 2009 [225]. A few months later, the same group reported the detection of the first recessive disorder (Miller syndrome) detected by whole-exome sequencing (WES) [226]. These two papers marked a new era where NGS became the preferred tool for rare Mendelian disease gene identification. There are several excellent reviews that describe the exponential growth in disease gene identification that started in 2010 [227–229]. Up to 27 February 2014, the number of genes with phenotype-causing mutations has reached 3162 according to online Mendelian inheritance in man (OMIM) Mgene map statistics [230]. In a recent review, Rabbani *et al.* estimated that from January 2010 to May 2012, over 100 causative genes in various Mendelian disorders have been identified by means of exome sequencing [231].

WES is now a valid and standard diagnostic approach for the identification of molecular defects in patients with suspected genetic disorders. This fact was demonstrated last year by a publication in the *New England Journal of Medicine* by the Medical Genetics Laboratory group of Baylor College of Medicine. The group reported the WES sequencing of 250 probands referred by physician, 98% of the cases were billed to the insurance. They reported a 25% molecular diagnostic rate (62 cases) [232]. In September 2013, the NIH funded four groups to explore the use of NGS for newborn screening [233]. With the cost per genome getting close to the US\$1000, it is becoming affordable to get sequenced at an early age, allowing for reanalysis of our genetic information at multiple intervals during the life of a person (Figure 3). A recent review outlines the approach, challenges, and benefits of such screening for adult genetic disease risks [2]. We also recently published a proof of concept project aimed to evaluate the benefits of screening healthy adults [234]. Our pilot project demonstrated that when WES is combined with medical and family history the findings are substantial. In a cohort of 81 unrelated individuals, we identified 271 recessive risk alleles (214 genes), 126 dominant risk alleles (101 genes) and three X-recessive risk alleles (three genes). In addition, we linked personal disease histories with causative disease genes in 18 volunteers [234].

SENS BIOLOGICZNY → MEDYCyna PERSONALIZOWANA



Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

nature medicine

[Explore content](#) ▾[About the journal](#) ▾[Publish with us](#) ▾

[nature](#) > [nature medicine](#) > [analyses](#) > article

Analysis | [Open access](#) | Published: 11 January 2024

Insights for precision oncology from the integration of genomic and clinical data of 13,880 tumors from the 100,000 Genomes Cancer Programme

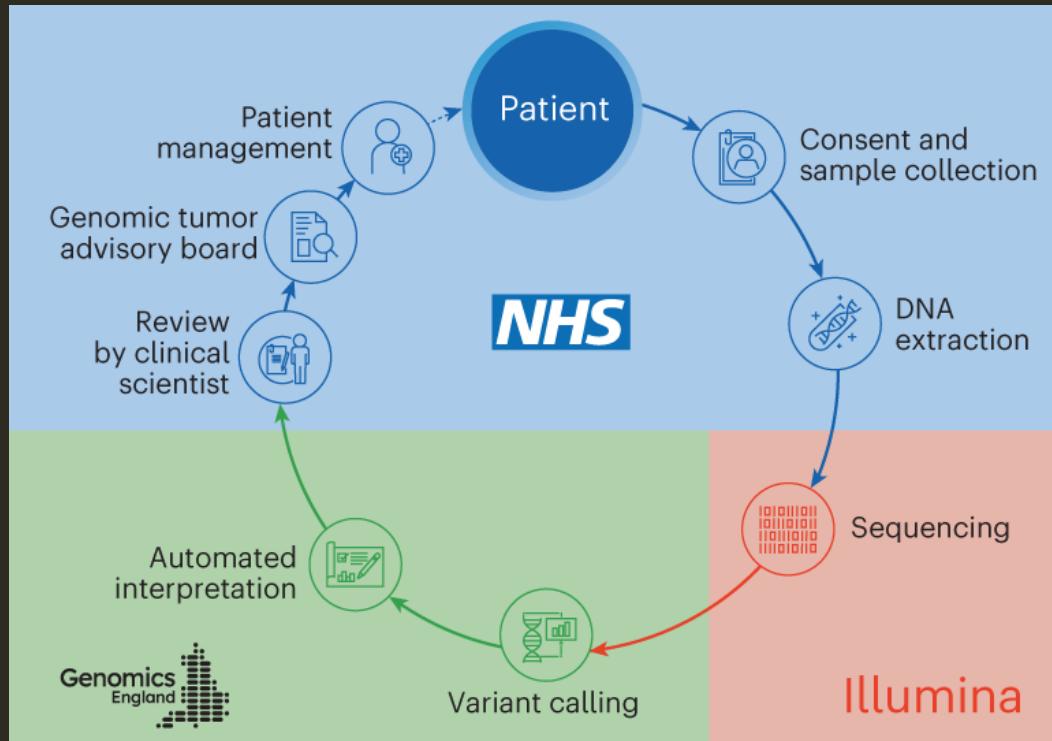
[Alona Sosinsky](#), [John Ambrose](#), [William Cross](#), [Clare Turnbull](#), [Shirley Henderson](#), [Louise Jones](#), [Angela Hamblin](#), [Prabhu Arumugam](#), [Georgia Chan](#), [Daniel Chubb](#), [Boris Noyvert](#), [Jonathan Mitchell](#), [Susan Walker](#), [Katy Bowman](#), [Dorota Pasko](#), [Marianna Buongermino Pereira](#), [Nadezda Volkova](#), [Antonio Rueda-Martin](#), [Daniel Perez-Gil](#), [Javier Lopez](#), [John Pullinger](#), [Afshan Siddiq](#), [Tala Zainy](#), [Tasnim Choudhury](#),
... [Nirupa Murugaesu](#)  + Show authors

The 100,000 Genomes Project, a UK Government initiative conducted within the National Health Service (NHS) in England, aimed to establish **standardized high-throughput whole-genome sequencing** (WGS) for patients with cancer and rare diseases via an automated, International Organization for Standardization-accredited bioinformatics pipeline (providing clinically accredited variant calling and variant prioritization).

Genomics England, alongside NHS England, analyzed WGS data from **13,880 solid tumors** spanning **33 cancer types**, integrating genomic data with real-world treatment and outcome data, within a secure Research Environment.

A longer-term objective was to accelerate the delivery of molecular testing, including WGS, in NHS clinical cancer care.

JOURNEY OF THE PATIENT'S GENOME



Patients provided written informed consent for WGS analysis.

DNA was extracted from **tumor** and **normal** (blood) samples using standardized protocols and samples were submitted for WGS, which was performed on an **Illumina** sequencer.

An automated **pipeline** was constructed for sequence quality control, alignment, variant calling and interpretation, with results returned to the 13 NHS Genomic Medicine Centers for review.

WYNIKI

Incidence of somatic mutations in genes recommended for standard-of-care testing varied across cancer types

- in glioblastoma multiforme, small variants were present in 94% of cases
- sarcoma demonstrated the highest occurrence of actionable structural variants (13%).
- Homologous recombination deficiency was identified in 40% of high-grade serous ovarian cancer cases with 30% linked to pathogenic germline variants, highlighting the value of combined somatic and germline analysis.

The linkage of WGS and longitudinal life course clinical data allowed the assessment of treatment outcomes for patients stratified according to pangenomic markers.

*glejak wielopostaciowy, mięsak ,rak jajnika

SENS BIOLOGICZNY → GENETIC COUNSELING

www.youtube.com/watch?v=C34dhHUWp0Y

A YouTube video player showing a woman with shoulder-length brown hair, wearing a red blazer over a white top, speaking directly to the camera. She has a slight smile and is looking slightly to her left. The background is a plain, dark grey. In the bottom right corner of the video frame, there is text identifying her: "Katherine Schneider, MPH, LGC" and "Senior Genetic Counselor". Below this, the logo for "DANA-FARBER CANCER INSTITUTE" is displayed, featuring a stylized 'F' icon followed by the text "DANA-FARBER" and "CANCER INSTITUTE". The video player interface includes a progress bar at the bottom left showing "0:15 / 5:09", and a standard set of control icons (play, pause, volume, etc.) at the bottom right.

Genetic Counseling for Cancer Risk: What to Expect | Dana-Farber Cancer Institute

Dana-Farber Cancer Institute
46,9 tys. subskrybentów

Subskrybuj

Like 10 Dislike Share Download Clip Save More

Best practices for variant calling in clinical sequencing

[Daniel C. Koboldt](#) 

[Genome Medicine](#) 12, Article number: 91 (2020) | [Cite this article](#)

148k Accesses | 134 Citations | 19 Altmetric | [Metrics](#)

Abstract

Next-generation sequencing technologies have enabled a dramatic expansion of clinical genetic testing both for inherited conditions and diseases such as cancer. Accurate variant calling in NGS data is a critical step upon which virtually all downstream analysis and interpretation processes rely. Just as NGS technologies have evolved considerably over the past 10 years, so too have the software tools and approaches for detecting sequence variants in clinical samples. In this review, I discuss the current best practices for variant calling in clinical sequencing studies, with a particular emphasis on trio sequencing for inherited disorders and somatic mutation detection in cancer patients. I describe the relative strengths and weaknesses of panel, exome, and whole-genome sequencing for variant detection. Recommended tools and strategies for calling variants of different classes are also provided, along with guidance on variant review, validation, and benchmarking to ensure optimal performance. Although NGS technologies are continually evolving, and new capabilities (such as long-read single-molecule sequencing) are emerging, the “best practice” principles in this review should be relevant to clinical variant calling in the long term.

Emphasis on **trio** sequencing for inherited disorders and somatic mutation detection in cancer patients

Strengths and weaknesses of:

- panel
- exome
- whole-genome sequencing for variant detection.

Recommended **tools** and **strategies** for calling variants of different classes

The “best practice” principles in this review should be relevant to clinical variant calling in the long term.

Best practices for variant calling in clinical sequencing

[Daniel C. Koboldt](#) 

[Genome Medicine](#) 12, Article number: 91 (2020) | [Cite this article](#)

NGS technologies enabled ambitious large-scale genomic sequencing efforts that have transformed our understanding of human health and disease, such as The Cancer Genome Atlas, the Centers for Mendelian Genomics.

They have also been widely adopted for *clinical genetic testing*:

- Targeted panels to interrogate medically relevant subsets of genes have become core components of precision oncology.
- Whole-exome sequencing, which selectively targets the protein-coding regions of known genes, has become a frontline diagnostic tool for inherited disorders

SENS BIOLOGICZNY → BADANIA PRENATALNE

OPIS TESTU

jest testem przesiewowym, nie jest testem diagnostycznym. Jego działanie polega na izolowaniu cfDNA (w tym zarówno matczynego, jak i płodowego DNA) z próbki krwi matki i wykonywaniu sekwencjonowania całego genomu przy użyciu technologii sekwencjonowania nowej generacji. Unikalne odczyty każdego chromosomu są obliczane i porównywane z optymalną próbką kontrolną stanowiącą punkt odniesienia. Dane są analizowane przy użyciu zastrzeżonych algorytmów bioinformatycznych BGI i ocena jest sporządzana tylko dla badanych warunków. Badania powinny być zawsze zamawiane przez wykwalifikowanego pracownika służby zdrowia a wyniki omawiane z pacjentem. Badanie nie może być stosowane jako jedyna podstawa do diagnozy lub innej decyzji w sprawie leczenia ciąży. Identyfikacja płci polega na izolowaniu DNA wolnego od komórek (w tym zarówno matczynego, jak i płodowego DNA) z próbki krwi matki, a następnie molekularnych badań genetycznych w celu określenia względnych ilości chromosomu Y.

ZASTRZEŻENIA

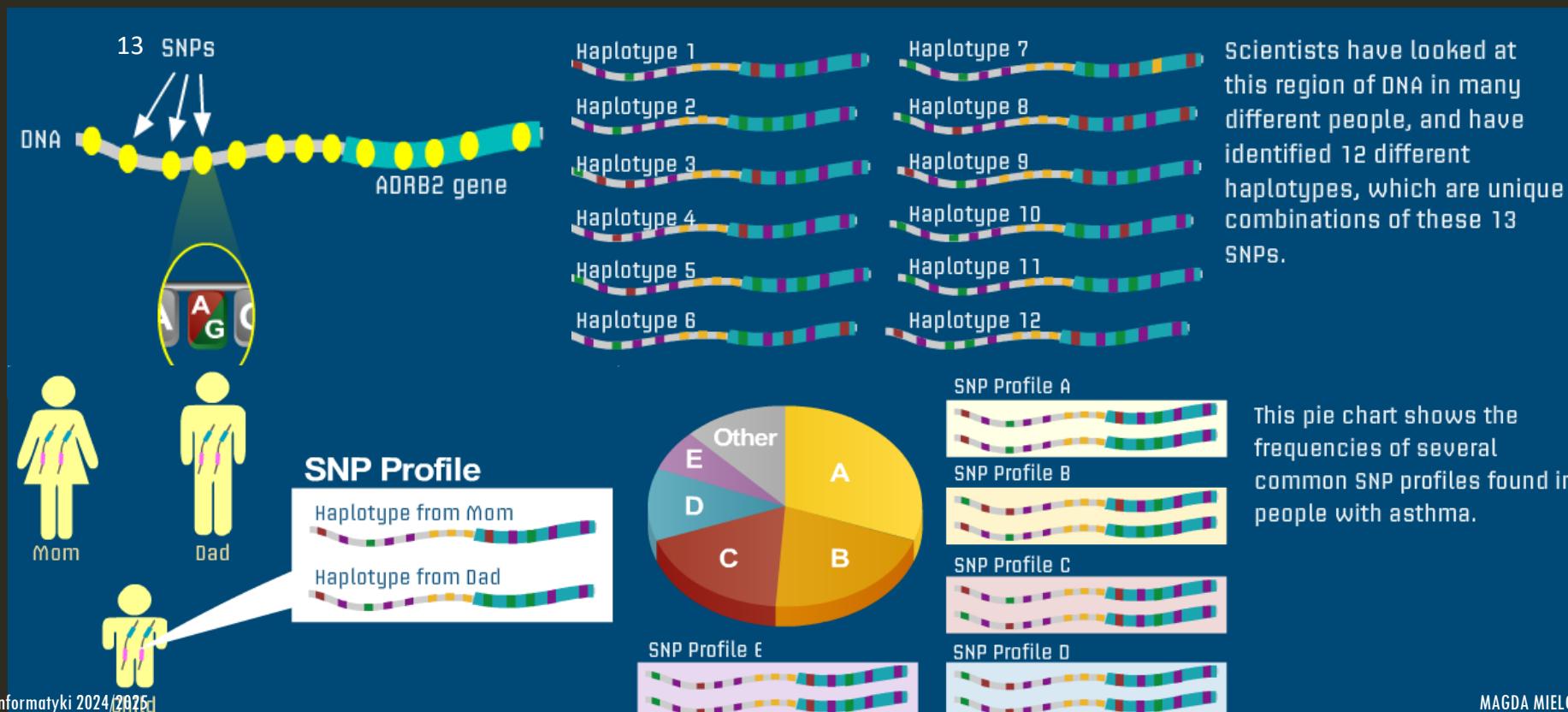
Test NIFTY Pro™ NIE jest testem diagnostycznym, wyniki są przeznaczone do wykorzystania informacyjnego, dlatego nie można wykluczyć wyników fałszywie dodatnich i fałszywie ujemnych. Wydajność "Inne ustalenia" nie została w pełni zweryfikowana, ale dane w poniższej tabeli mogą być wykorzystane do celów informacyjnych. W tym teście wykryto 84 rodzaje zespołów del/dup. Niektóre z chorób na liście zespołów del/dup mogą być również spowodowane przez inne czynniki genetyczne, NIFTY ProTM tylko wykrywa i analizuje konkretny fragment zgodnie z autoryzowanymi bazami danych. Wykrywanie chromosomu Y zawarte w niniejszym raporcie nie może być wykorzystywane do diagnozowania płci płodowej lub chorób związanych z płcią i jest wykorzystywane wyłącznie jako dodatkowe informacje do analizy referencyjnej. Potencjalne źródła niedokładnego wyniku badania mogą obejmować między innymi: mozaicyzm matki, płodu i/lub łożyska, niską frakcję płodową, transfuzję krwi, operację przeszczepu, terapię komórkami macierzystymi, terapię heparyną i nieprawidłowy kariotyp biologicznych rodziców lub surrogatki. Wynik badania jest specyficzny dla badanej próbki i powinien być zawsze interpretowany przez wykwalifikowanego specjalistę w kontekście danych klinicznych i rodzinnych.

ZABURZENIE	CZUŁOŚĆ
Trisomia 21	99,17%
Trisomia 18	98,24%
Trisomia 13	>99,9%
Trisomia 9, 16, 22	NA
CNV	>90%
Płeć	99,53%
ANEUPLOIDIE CHROMOSOMÓW PŁCIOWYCH	CZUŁOŚĆ
XYY	>99,9%
XXY	>99,9%
XXX	>99,9%
XO	>99,9%

SENS BIOLOGICZNY → FARMAKOGENETYKA/FARMAKOGENOMIKA

„Applying SNP profiles to drug choices”

- astma ← gen ADRB2
- polimorfizmy punktowe a odpowiedź na leczenie albuterolem
- albuterol łagodzi objawy astmy tylko u niektórych chorych



SENS BIOLOGICZNY → FARMAKOGENETYKA/FARMAKOGENOMIKA

The diagram illustrates the relationship between SNP profiles and drug response, leading to personalized medicine.

SNP Profile Distribution:

SNP Profile	Albuterol Response
A	Poor
B	Good
C	Fair
D	None
E	Very Good

LAB REPORT (Patient Johnson):

- Patient Johnson
- SNP Profile: D
- NOT RESPONSIVE TO ALBUTEROL

Future Scenario:

SNP Profile	Albuterol Response
A	Poor
B	Good
C	Fair
D	None
E	Very Good

In the future, a physician will be able to determine a patient's SNP profile, compare it with known data, and predict whether the patient will respond to the drug albuterol.

Illustration: A red 'X' is drawn over two inhalers labeled "ALBUTEROL" and "Other Drug".

The physician can then design the patient's treatment accordingly. This will be a great improvement over the trial-and-error method physicians use today.

NGS TOOLS & PIPELINES

The screenshot shows the Galaxy web interface. On the left, a sidebar lists various tools and pipelines: Upload, Tools (selected), Workflows, Visualization, Histories, and Pages. The main area displays a search bar for 'All Tools' and a summary text about Galaxy. A modal window is open, showing Nextflow pipeline code:

```
process sayHello {  
    input:  
        val cheers  
    output:  
        stdout  
    ....  
    echo $cheers  
    ....  
}  
  
workflow {  
    channel.of('Ciao','Hello','Hola') | sayHello | video  
}
```

The modal also includes a small logo for 'GBCC2025 Cold Spring Harbor'.

Nextflow

Data-driven computational pipelines

Nextflow enables scalable and reproducible scientific workflows using software containers. It allows the adaptation of pipelines written in the most common scripting languages.

Its fluent DSL simplifies the implementation and the deployment of complex parallel and reactive workflows on clouds and clusters.

Pipelines



Browse the 104 pipelines that are currently available as part of nf-core.

FORMATY DANYCH - PODSUMOWANIE

