

Analiza danych z wysokoprzepustowego
sekwencjonowania

W4: genomy



sekwencjonowanie genomów

Zakończone projekty sekwencjonowania genomów

153 899

31 920

4 032

Genomes



111 840

26 450

3 672

Bacterial

1 509

649

177

Archaeal

40 550

4 821

183

Eukaryal

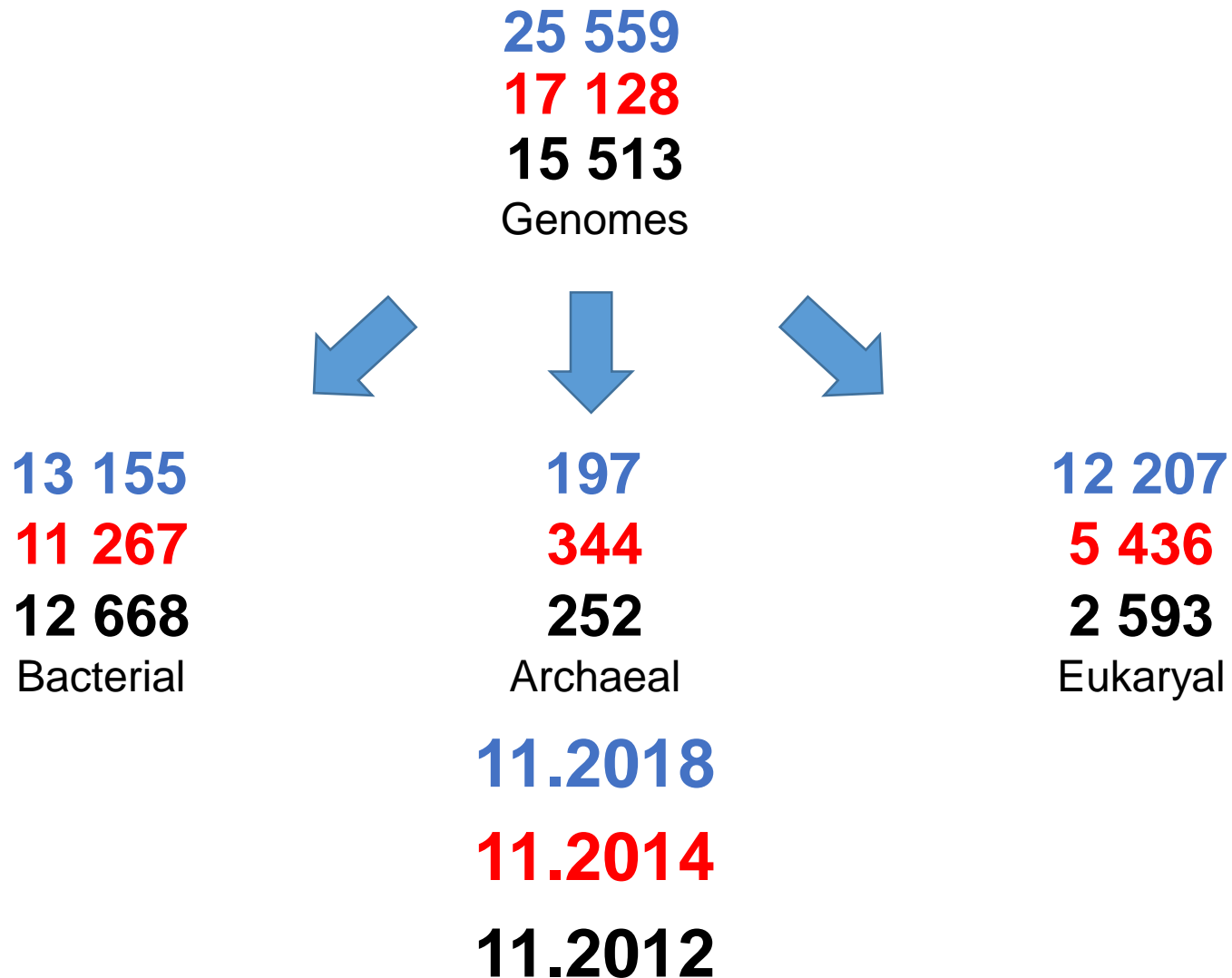
11.2018

11.2014

11.2012

sekwencjonowanie genomów

Trwające projekty sekwencjonowania genomów



sekwencjonowanie genomów

Sekwencjonowanie nowych organizmów:

- poznanie genomów organizmów modelowych
- identyfikacja nowych mechanizmów adaptacji do środowiska
- śledzenie ewolucji genomów
- inżynieria genetyczna organizmów istotnych gospodarczo i przemysłowo
- bo możemy i nas stać!

Sekwencjonowanie genomów ludzkich:

- poznanie różnic genetycznych na poziomie populacyjnym
- poznanie mechanizmów chorób genetycznych wielogenowych i sposobu ich dziedziczenia
- identyfikacja markerów chorobowych
- identyfikacja markerów prognostycznych w terapii

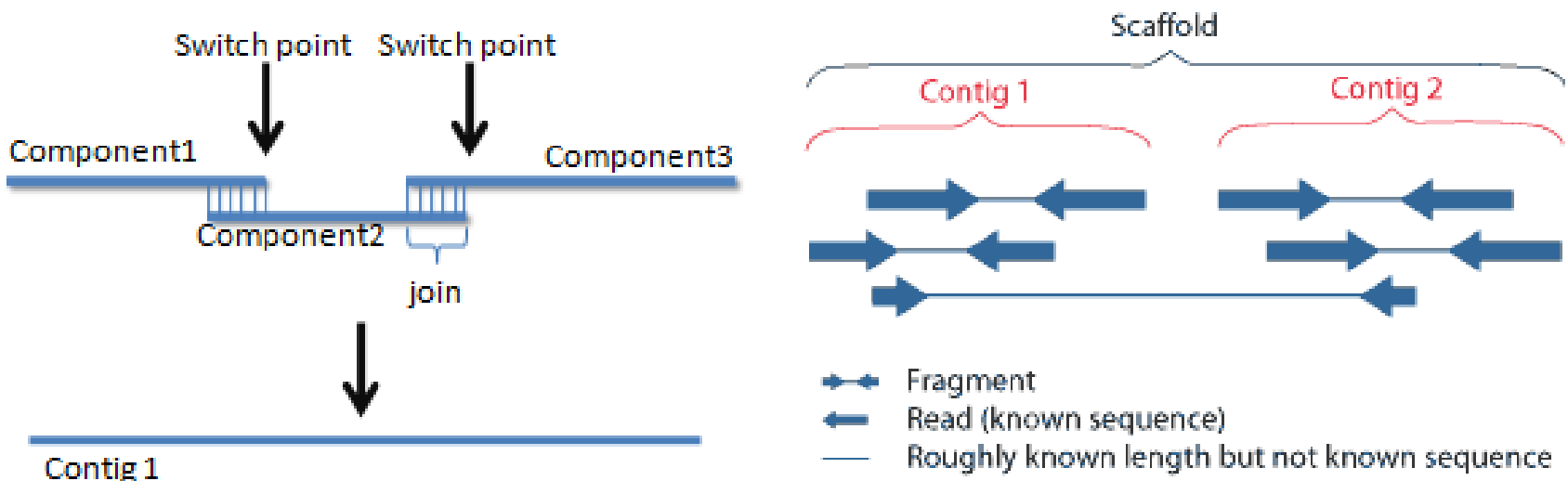
sekwencjonowanie genomów

terminologia

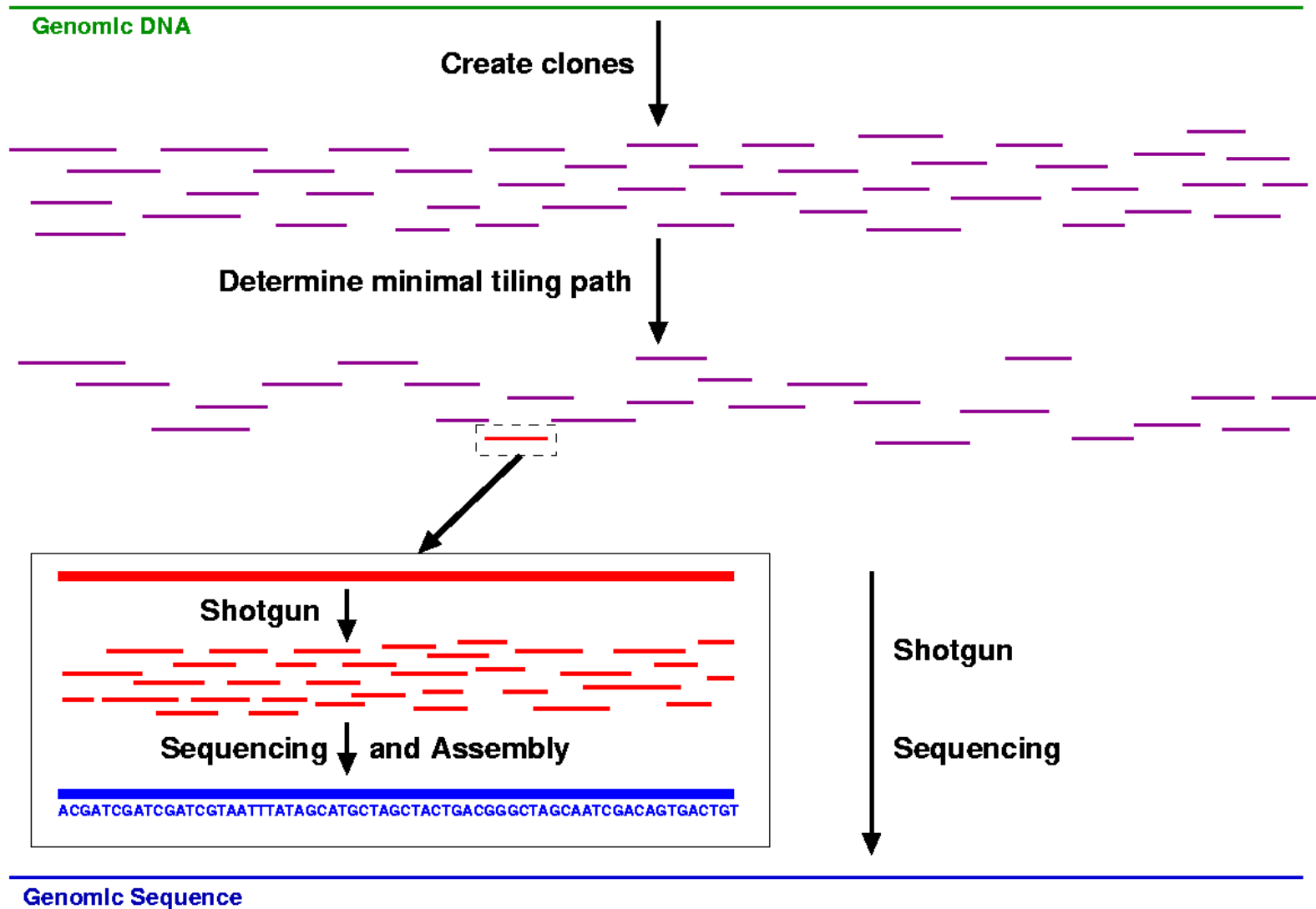
BAC – Bacterial Artificial Chromosome – plazmid zdolny do przenoszenia insertów o długości 150-350 kbp

contig – ciągła sekwencja wygenerowana na podstawie nakładających się na siebie fragmentów sekwencji uzyskanych z sekwencjonowania

scaffold – uporządkowany pod względem kolejności i orientacji zestaw contigów. Zazwyczaj zawiera przerwy w sekwencji.



sekwencjonowanie genomów hierarchiczne – Human Genome Project



sekwencjonowanie genomów wysokoprzepustowe

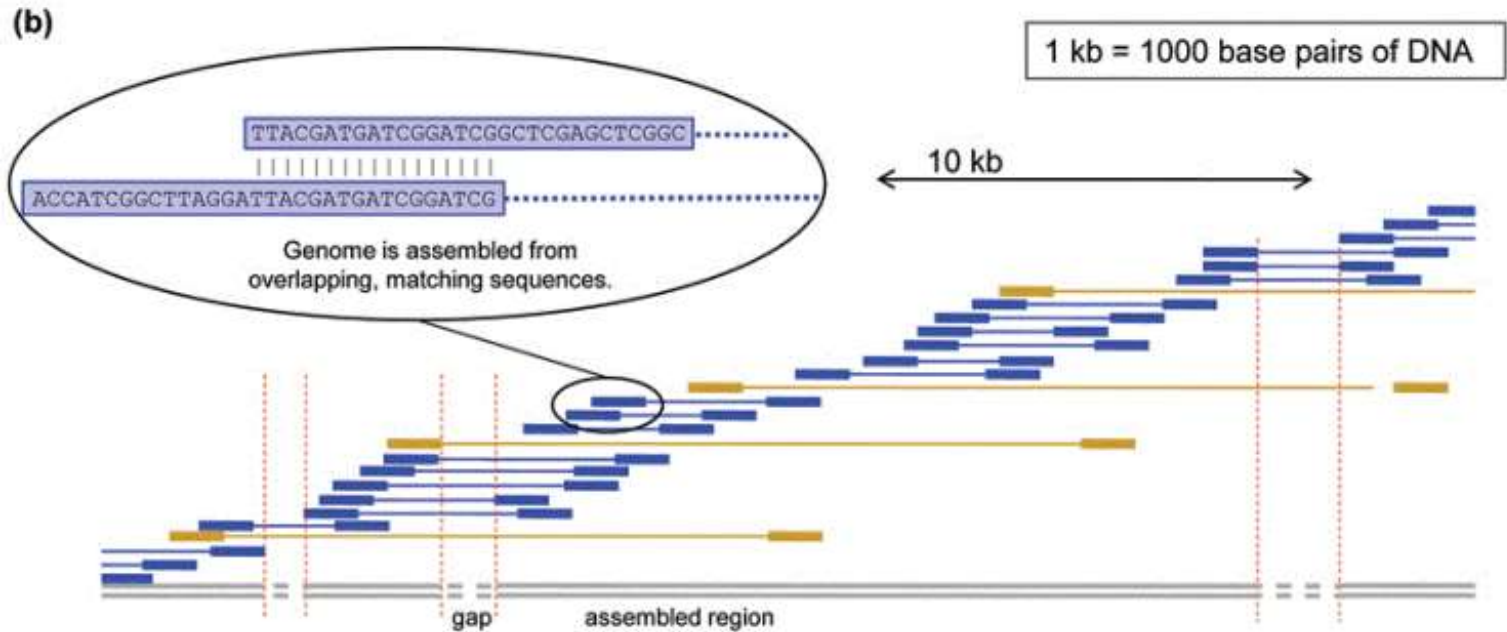
Zalety NGS:

- tanie
- ultra-wysokowydajne
- amplifikacja zamiast klonowania

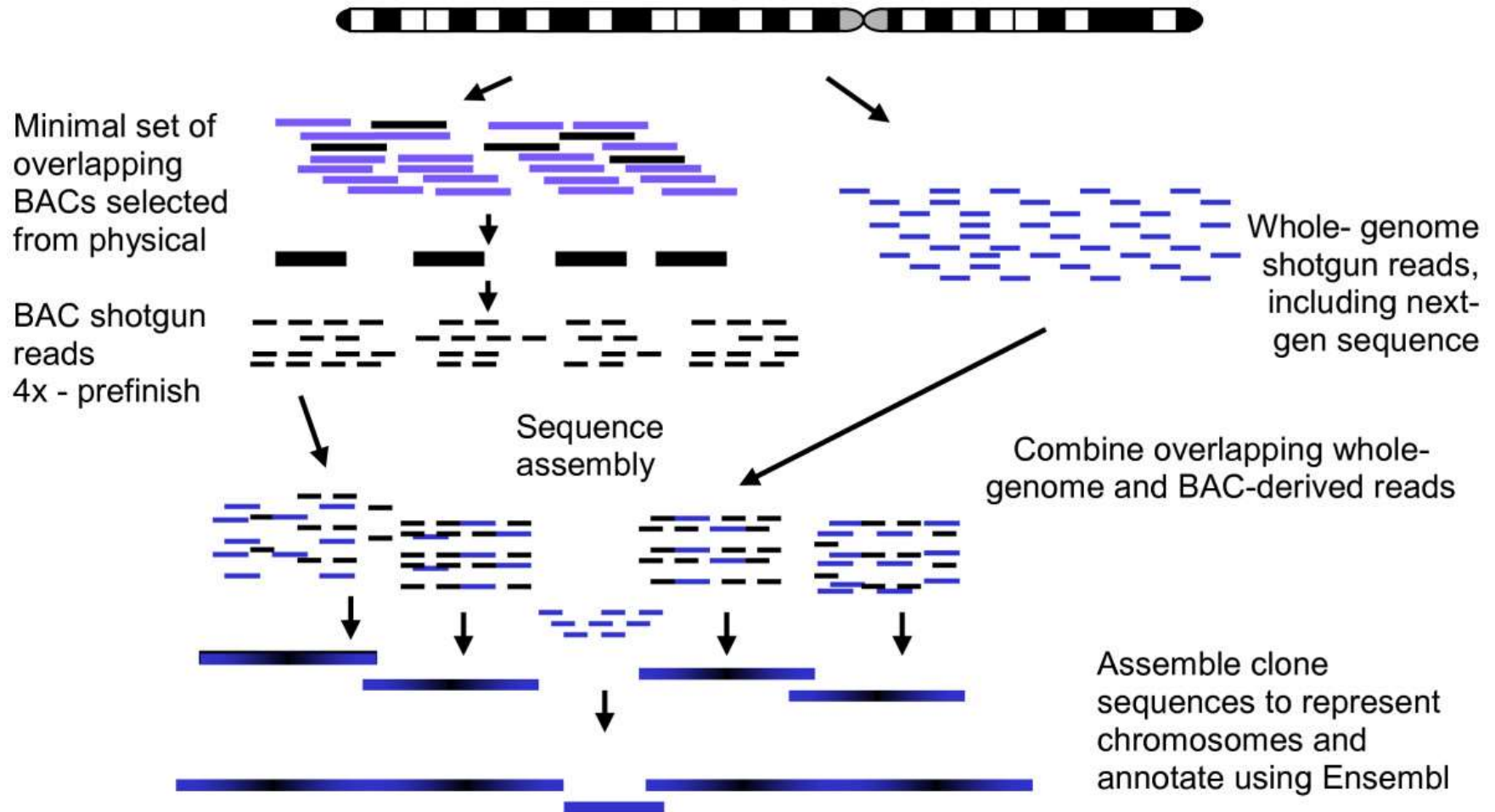
Wady:

- krótkie odczyty
- ograniczony dystans pomiędzy paired-ends

sekwencjonowanie genomów shotgun (NGS)



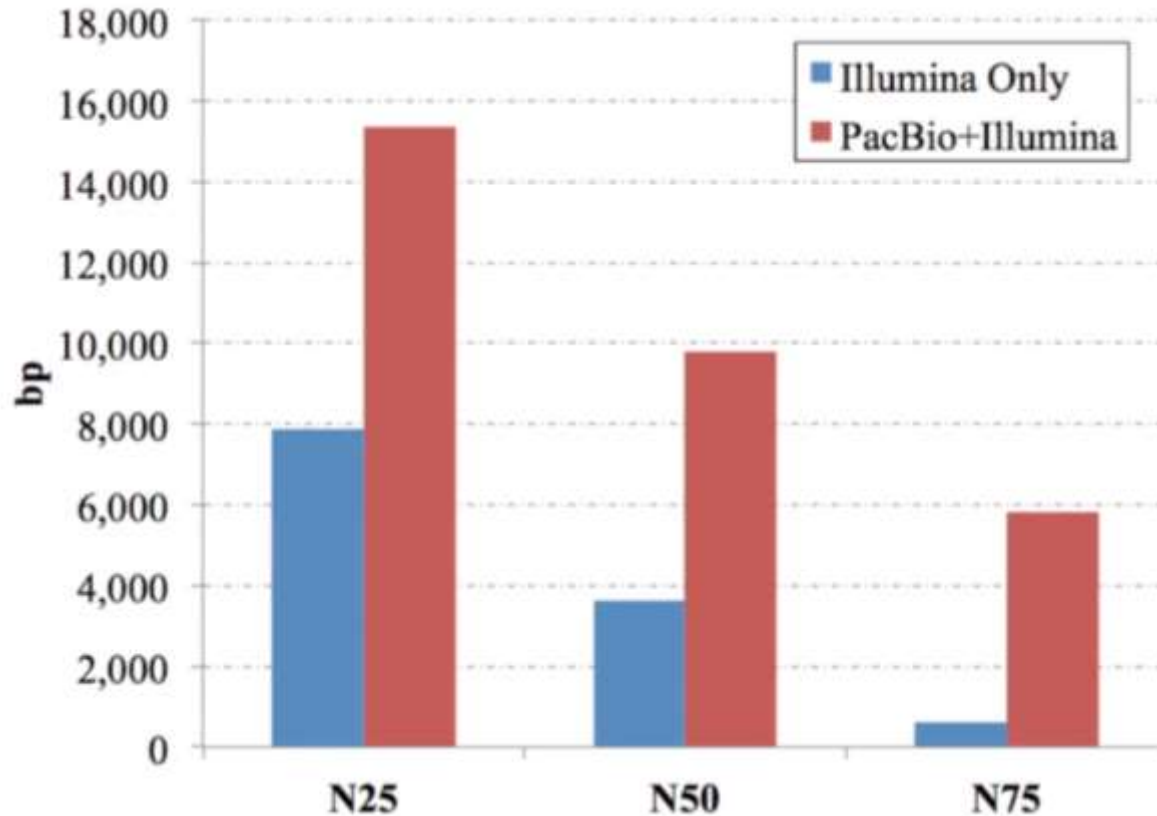
sekwencjonowanie genomów podejście hybrydowe



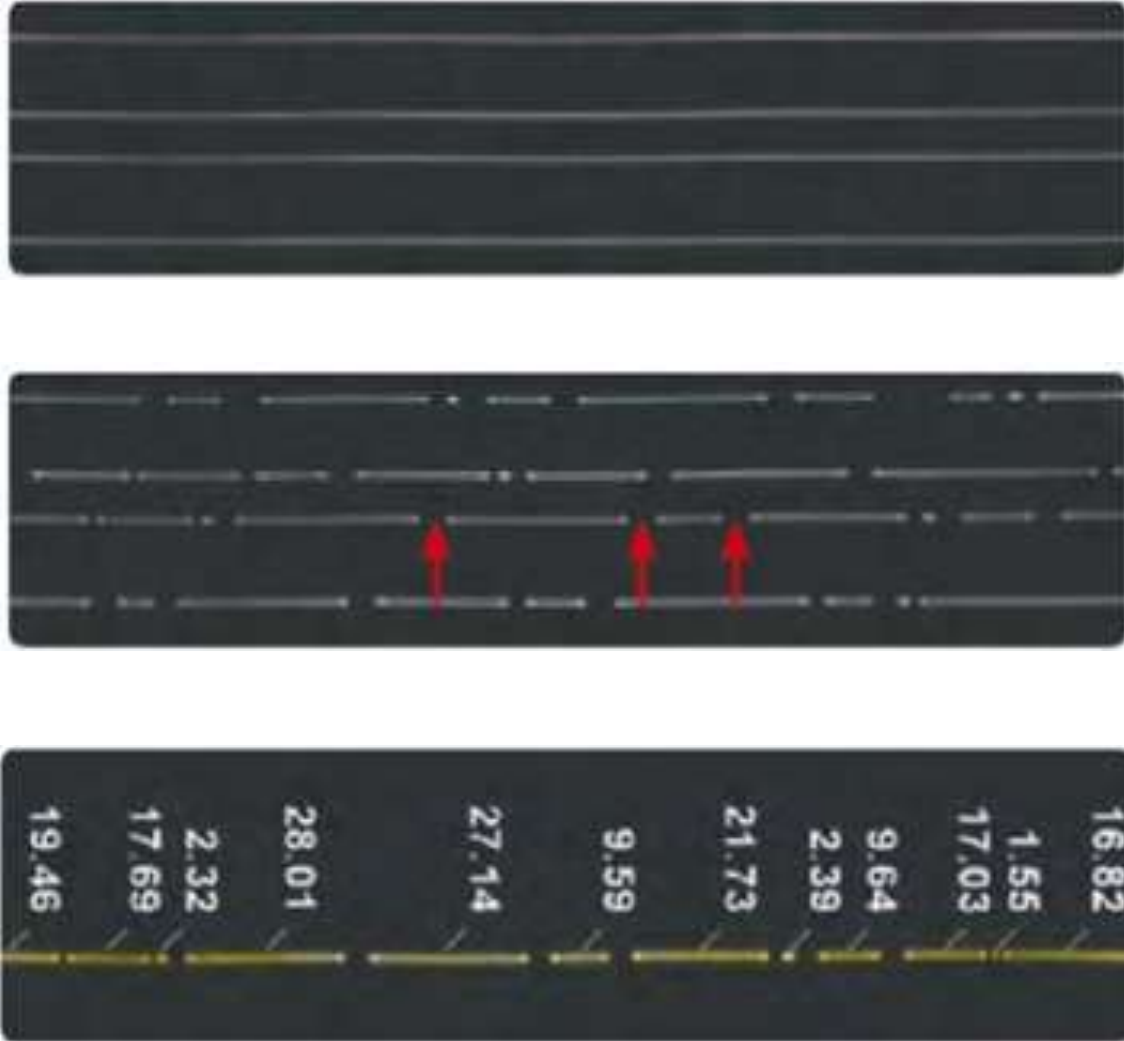
sekwencjonowanie genomów alternatywy dla BAC

- Wprowadzenie długich odczytów jako rusztowania:
 - 454
 - PacBio
 - Oxford Nanopore
- Mapowanie optyczne genomu

sekwencjonowanie genomów podejście hybrydowe PacBio

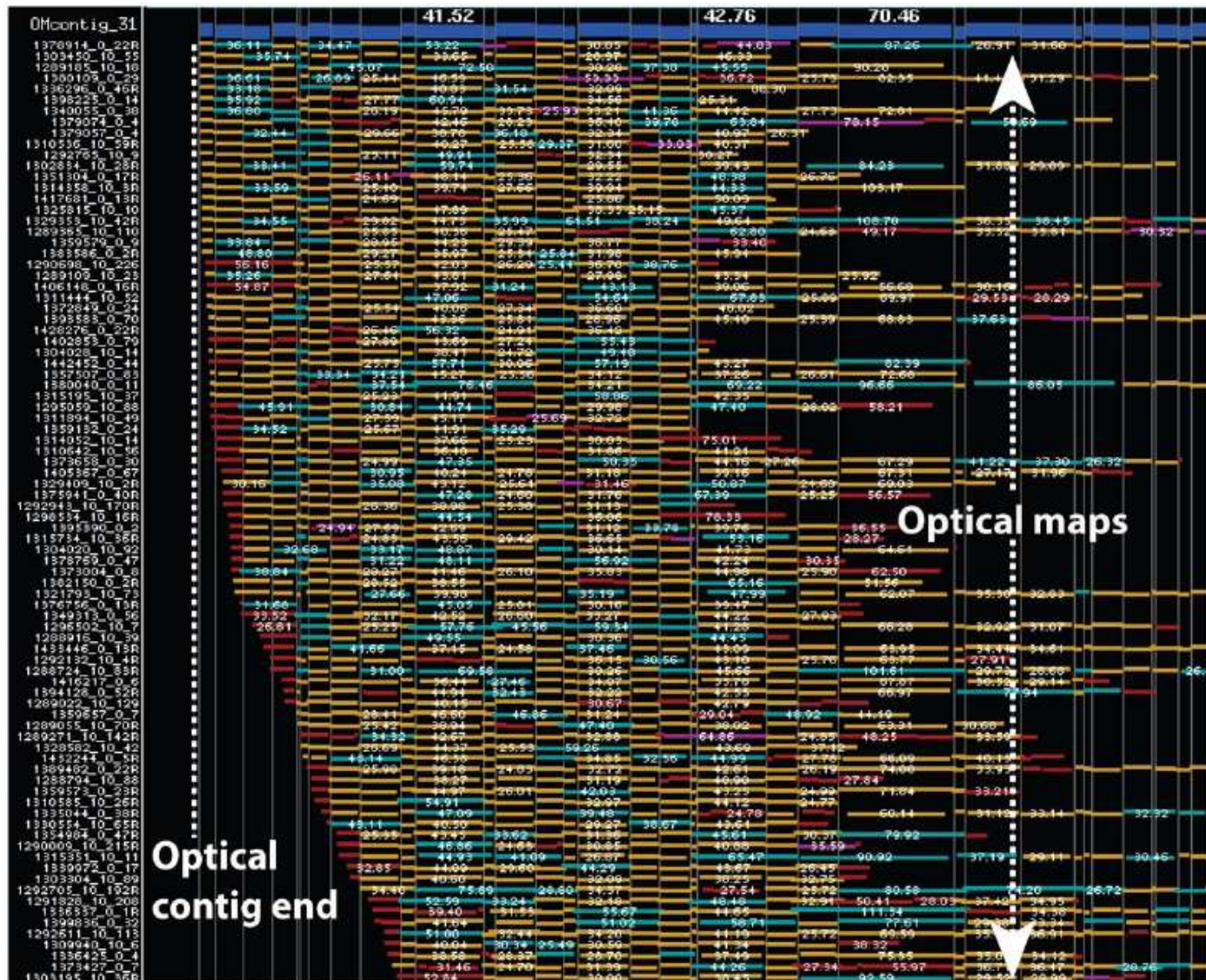


sekwencjonowanie genomów mapowanie optyczne



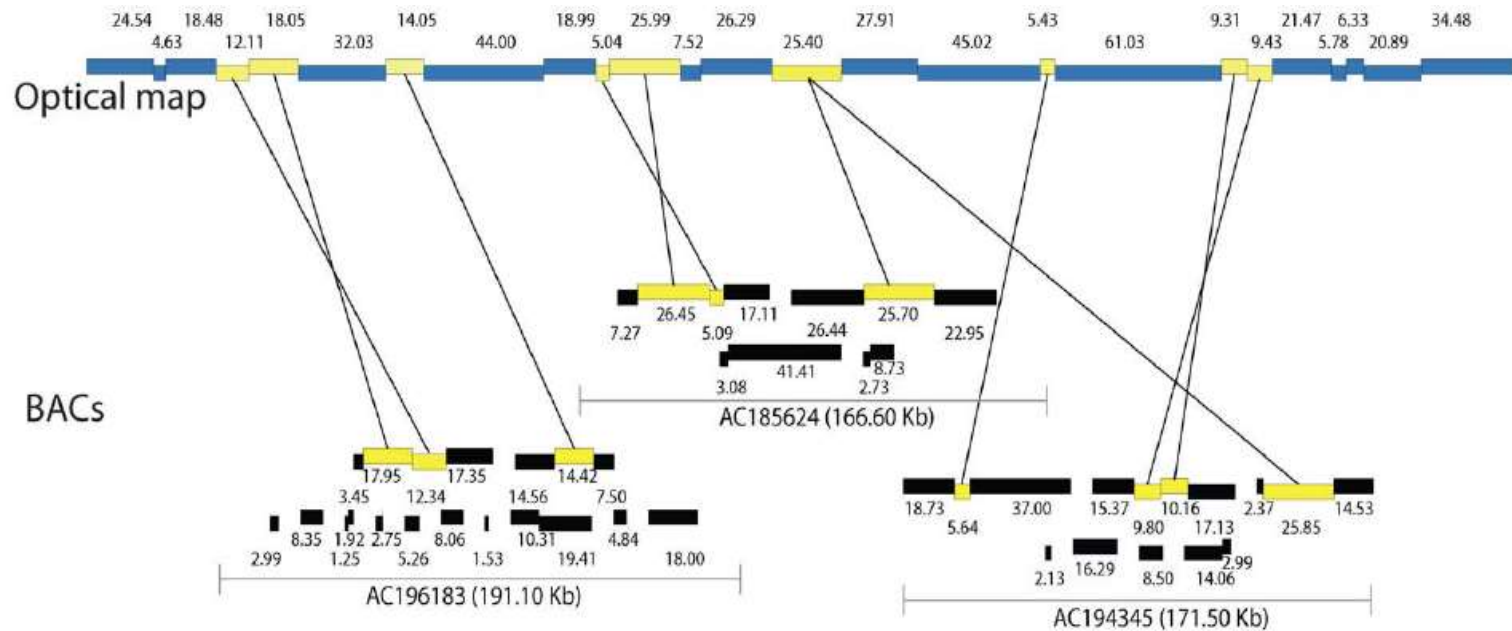
sekwencjonowanie genomów

mapowanie optyczne



sekwencjonowanie genomów

mapowanie optyczne



sekwencjonowanie genomów wysokoprzepustowe

shotgun

- genomy bakteryjne
- małe genomy eukariotyczne
- duże genomy eukariotyczne z niską ilością sekwencji powtórzonych

hierarchiczne/hybrydowe/wspomagane

- duże genomy eukariotyczne o znacznej liczbie sekwencji powtórzonych (ssaki, rośliny)

składanie sekwencji *de novo*

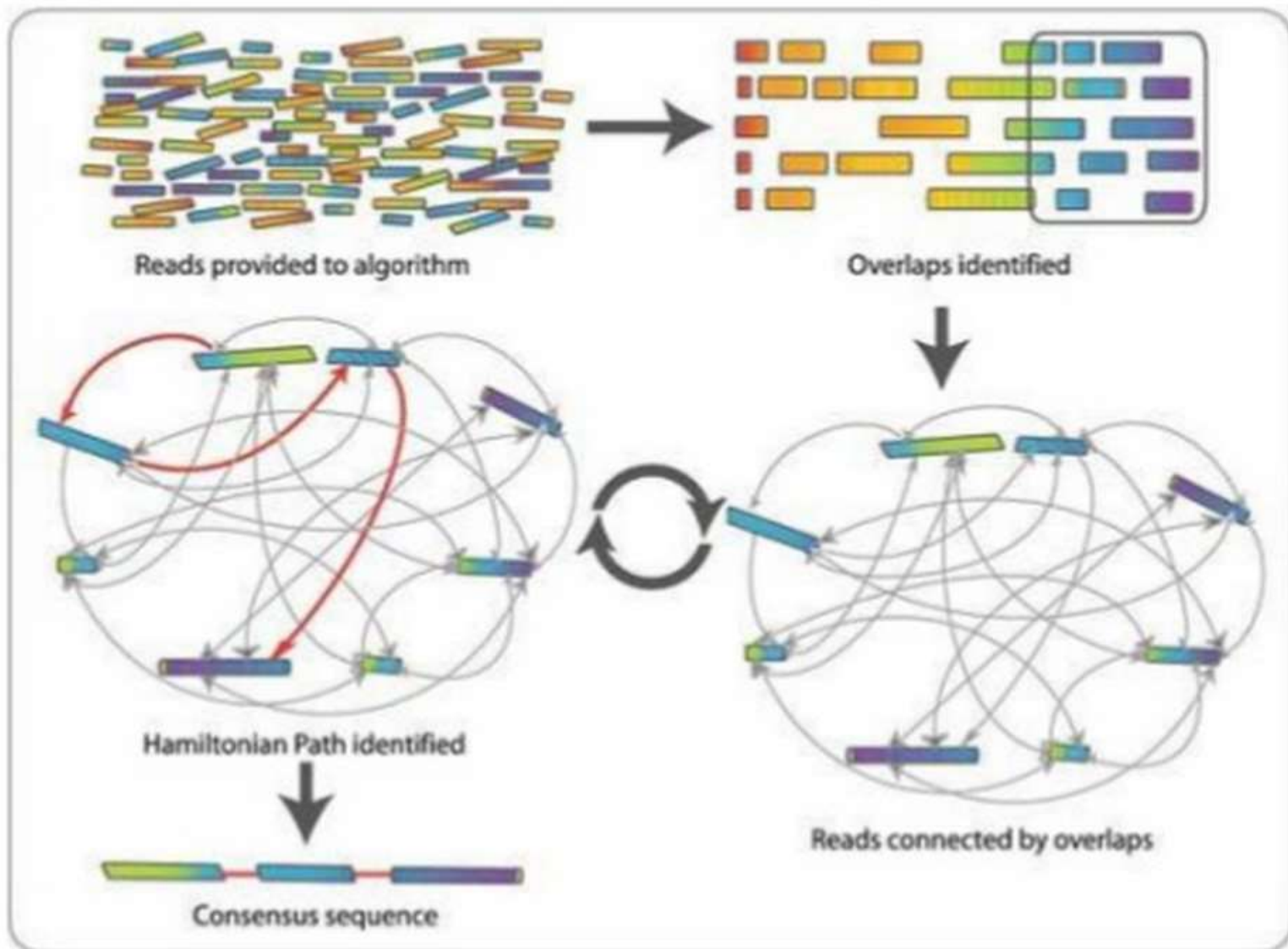
**proces rekonstrukcji oryginalnego DNA
z użyciem fragmentów sekwencji**

**Wszystkie metody robią to samo
Różnią się skrótem który biorą**

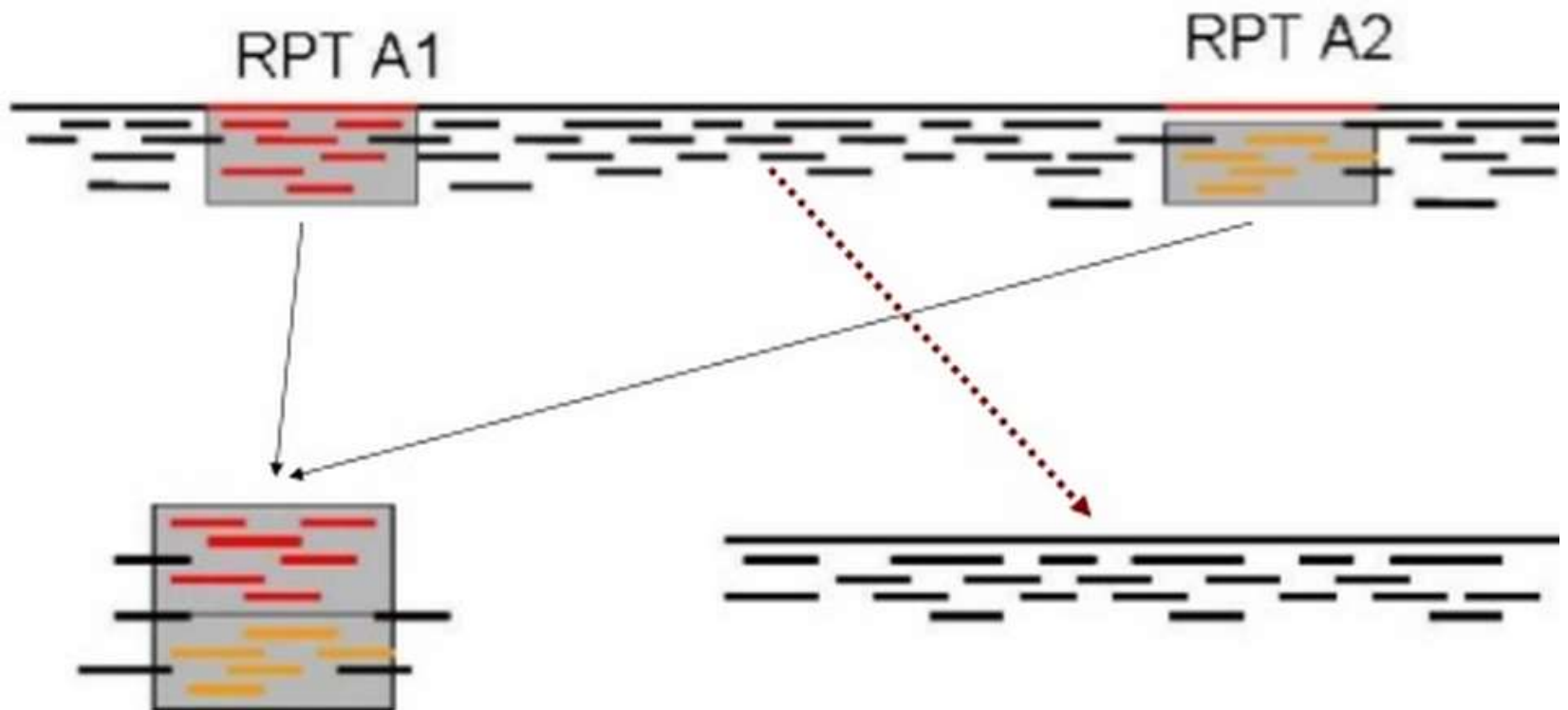
- składanie zachłanne
- poszukiwanie nałożenia sekwencji
- grafy de Bruijn'a
- grafy ciągów
- „seed and extend”

- Znaleźć wszystkie możliwe nałożenia pomiędzy odczytami
- Zbudować graf
- Uprościć graf
- Przejść graf

workflow



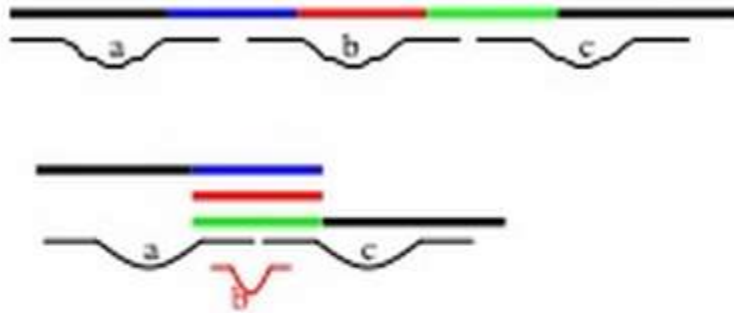
powtórzenia



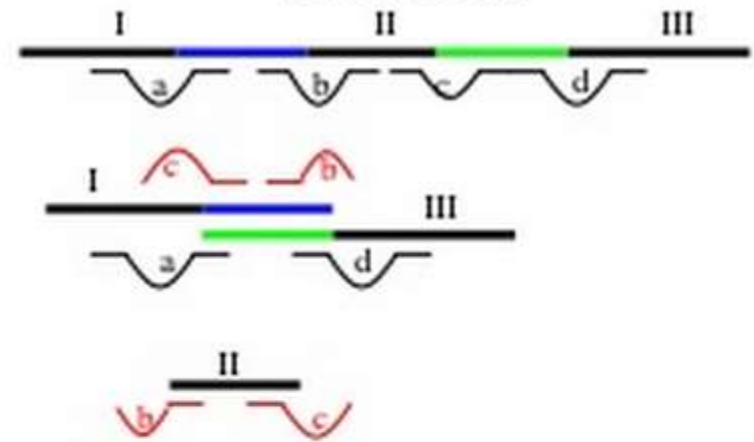
The repeated element is collapsed into a single contig

powtórzenia

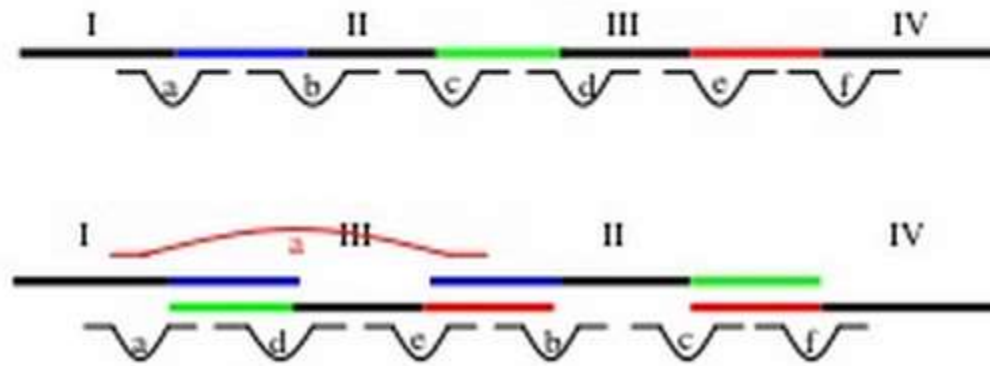
collapsed tandem



excision



rearrangement



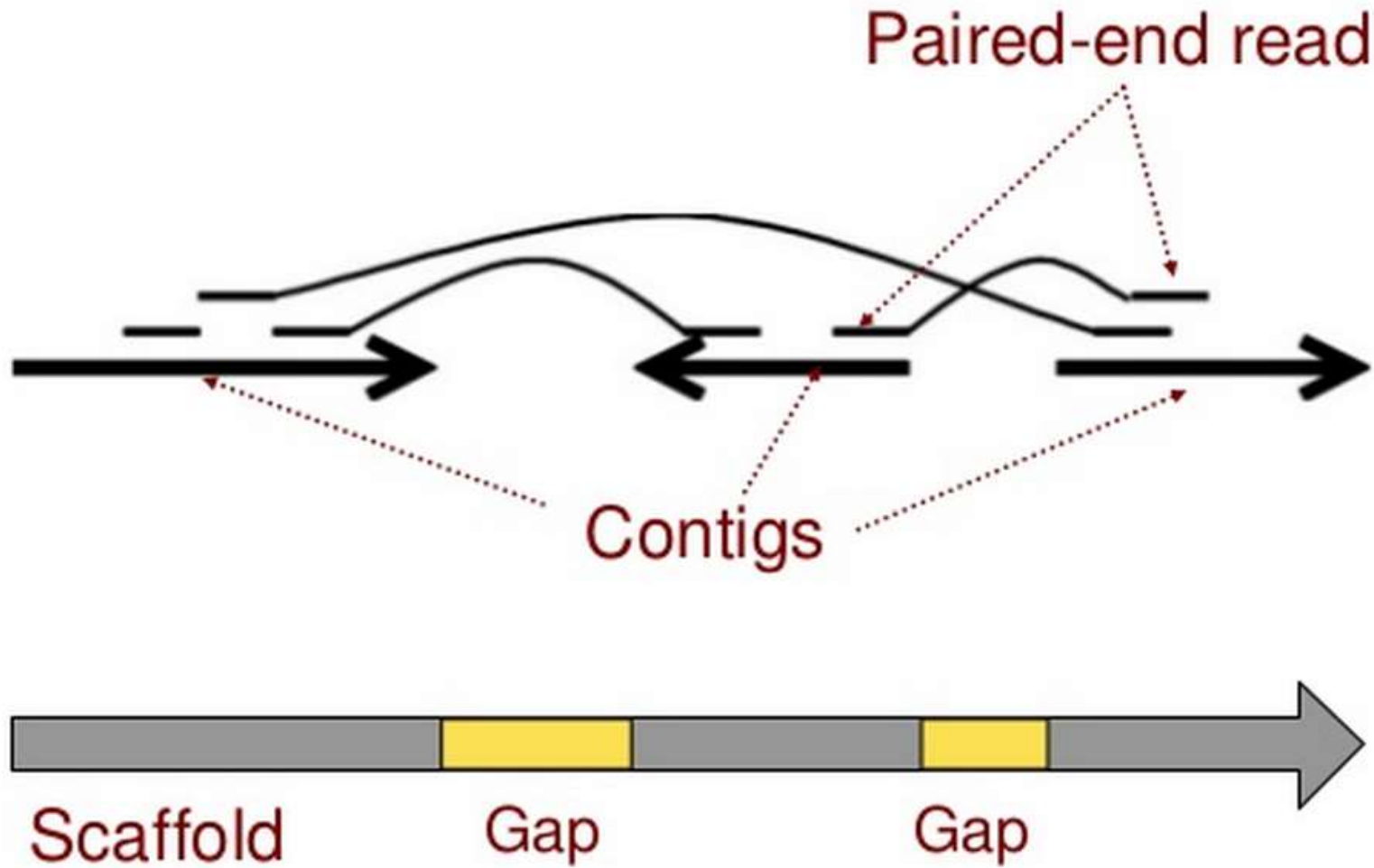
3 prawa powtórzeń

1. Niemożliwe jest rozwiązanie struktury powtórzenia o długości S dopóki nie posiadasz odczytów dłuższych niż S
2. Niemożliwe jest rozwiązanie struktury powtórzenia o długości S dopóki nie posiadasz odczytów dłuższych niż S
3. Niemożliwe jest rozwiązanie struktury powtórzenia o długości S dopóki nie posiadasz odczytów dłuższych niż S

odczyty paired-end

- bardzo pomocne podczas tworzenia scaffoldów
- mogą pochodzić z różnej długości fragmentów
- te z krótkich najczęściej są składane w ramach jednego kontigu – przyczyniają się do wydłużenia kontigu
- te z długich czasem dzielą się pomiędzy kontigi – jest to przesłanka do połączenie tych kontigów w scaffold

odczyty paired-end



miara sukcesu – N50

Przykład:

- mamy 10 kontigów o długościach: 1, 1, 3, 3, 5, 7, 8, 12, 16, 20
- Suma: 76 (tyle nukleotydów jest w naszym złożeniu)
- patrzymy „od góry” który kontig przekracza połowę sumy, czyli 38:

$$20+16 = 36 (<38)$$

$$20+16+12 = 48 (>38) \quad \mathbf{N50 = 12}$$

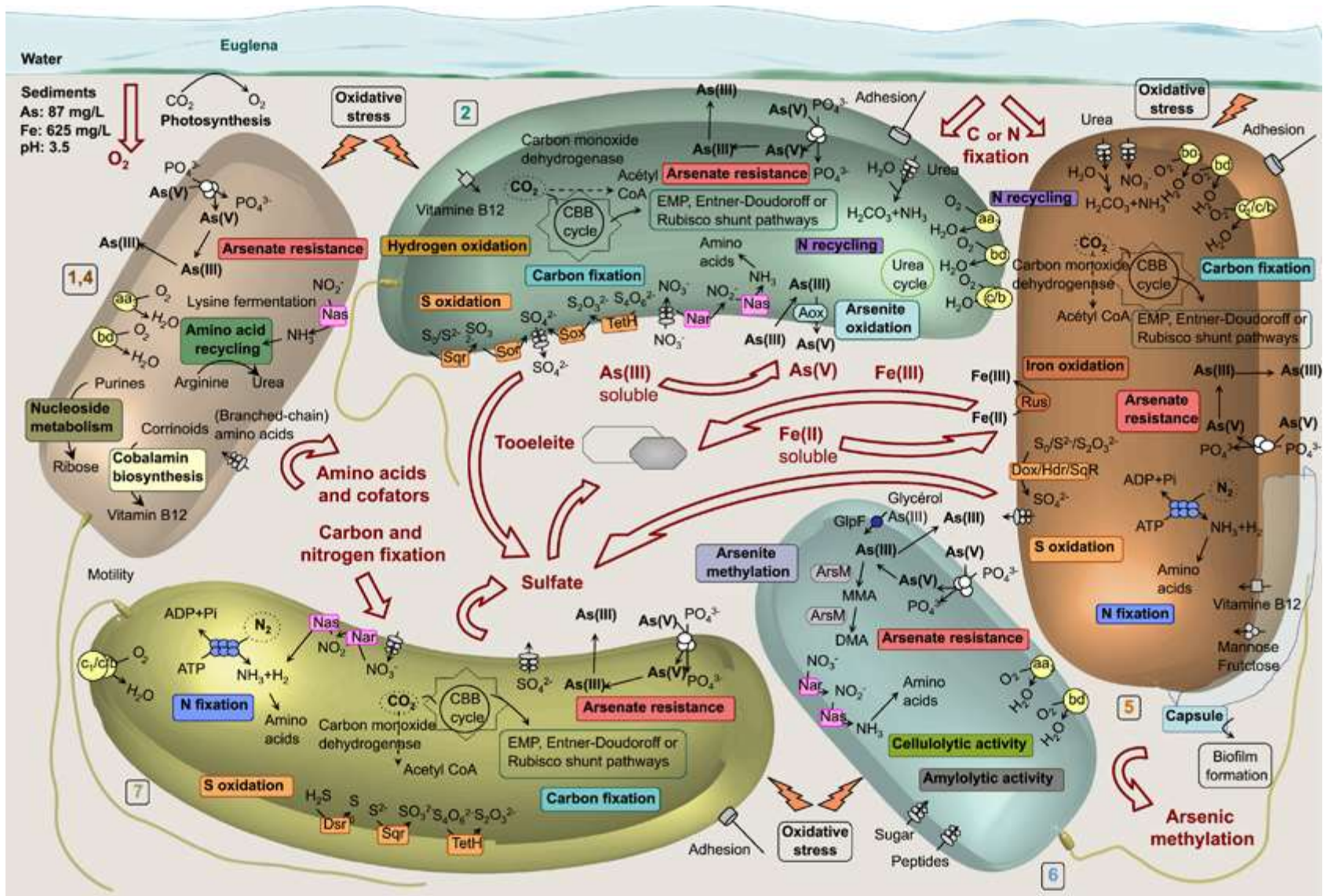
N50: długość kontigu, który wraz z krótszymi od niego zawiera połowę nukleotydów w złożeniu (połowę sumy długości wszystkich kontigów)

polecane narzędzia:

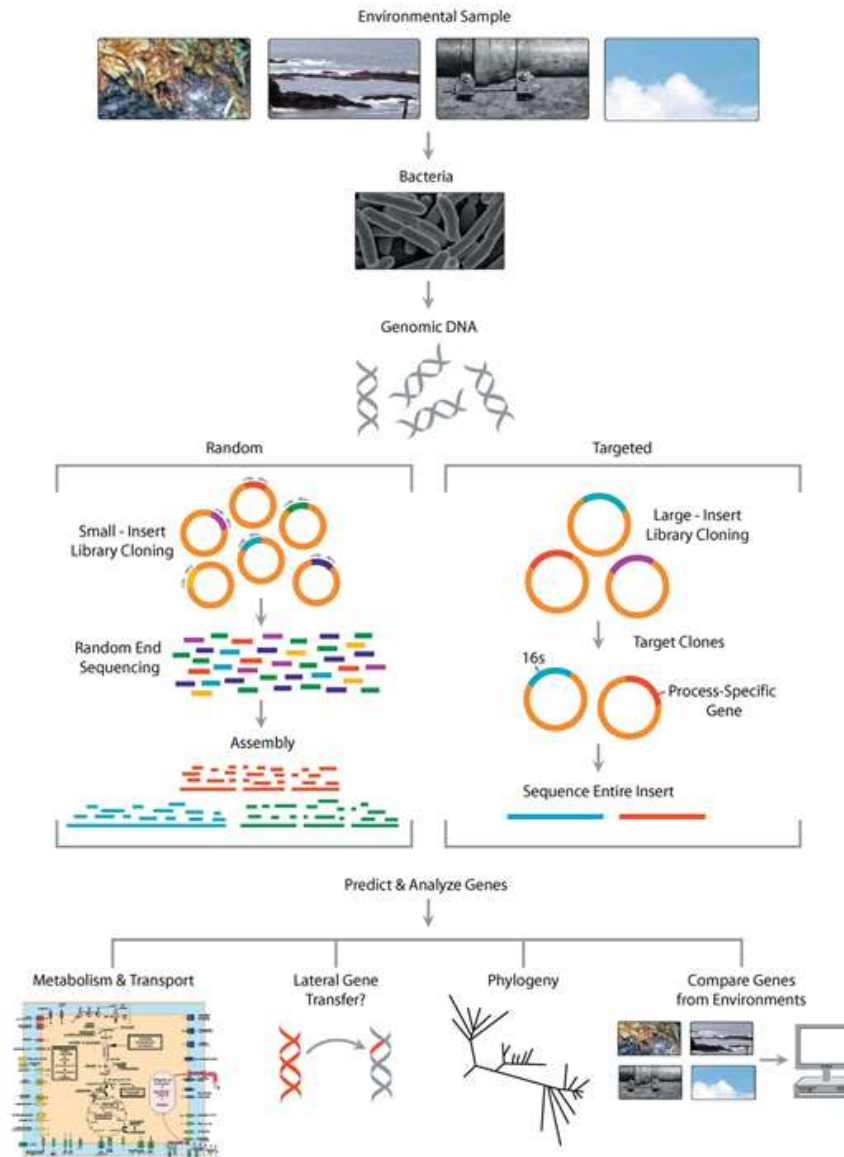
- Velvet
- Soap *de novo*
- aBYss
- Celera assembler
- Canu (PacBio, Nanopore)

sekwencjonowanie metagenomów

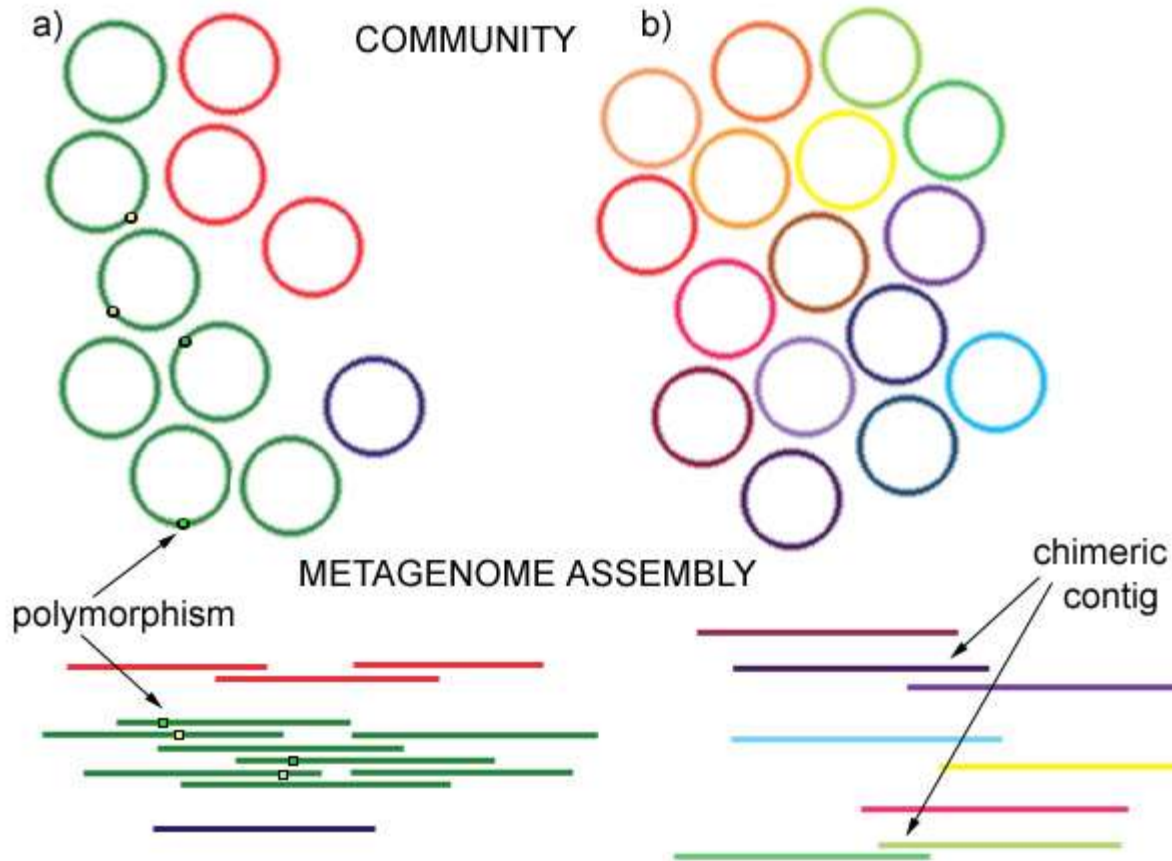
Metabolic diversity among main microorganisms inside an arsenic-rich ecosystem revealed by meta- and proteo-genomics



sekwencjonowanie metagenomów



sekwencjonowanie metagenomów



ARTICLES

A human gut microbial gene catalogue established by metagenomic sequencing

Junjie Qin^{1*}, Ruiqiang Li^{1*}, Jeroen Raes^{2,3}, Manimozhiyan Arumugam², Kristoffer Solvsten Burgdorf⁴, Chaysavanh Manichanh⁵, Trine Nielsen⁴, Nicolas Pons⁶, Florence Levenez⁶, Takuji Yamada², Daniel R. Mende², Junhua Li^{1,7}, Junming Xu¹, Shaochuan Li¹, Dongfang Li^{1,8}, Jianjun Cao¹, Bo Wang¹, Huiqing Liang¹, Huisong Zheng¹, Yinlong Xie^{1,7}, Julien Tap⁶, Patricia Lepage⁶, Marcelo Bertalan⁹, Jean-Michel Batto⁶, Torben Hansen⁴, Denis Le Paslier¹⁰, Allan Linneberg¹¹, H. Bjørn Nielsen⁹, Eric Pelletier¹⁰, Pierre Renault⁶, Thomas Sicheritz-Ponten⁹, Keith Turner¹², Hongmei Zhu¹, Chang Yu¹, Shengting Li¹, Min Jian¹, Yan Zhou¹, Yingrui Li¹, Xiuqing Zhang¹, Songgang Li¹, Nan Qin¹, Huanming Yang¹, Jian Wang¹, Søren Brunak⁹, Joel Doré⁶, Francisco Guarner⁵, Karsten Kristiansen¹³, Oluf Pedersen^{4,14}, Julian Parkhill¹², Jean Weissenbach¹⁰, MetaHIT Consortium†, Peer Bork², S. Dusko Ehrlich⁶ & Jun Wang^{1,13}

To understand the impact of gut microbes on human health and well-being it is crucial to assess their genetic potential. Here we describe the Illumina-based metagenomic sequencing, assembly and characterization of **3.3 million** non-redundant microbial genes, derived from 576.7 gigabases of sequence, from faecal samples of **124 European** individuals. The gene set, **~150 times larger than the human gene complement**, contains an overwhelming majority of the prevalent (more frequent) microbial genes of the cohort and probably includes a large proportion of the prevalent human intestinal microbial genes. The genes are largely shared among individuals of the cohort. Over 99% of the genes are bacterial, indicating that the entire cohort harbours between **1,000 and 1,150 prevalent bacterial species** and each individual at least 160 such species, which are also largely shared. We define and describe the minimal gut metagenome and the minimal gut bacterial genome in terms of functions present in all individuals and most bacteria, respectively.

sekwencjonowanie metagenomów

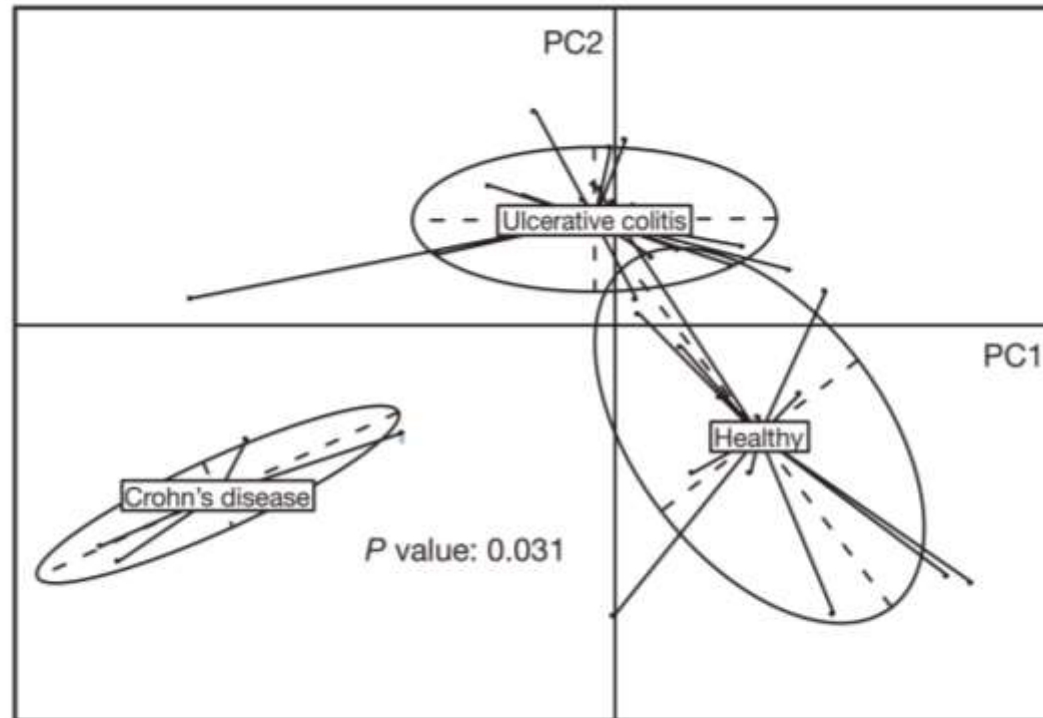


Figure 4 | Bacterial species abundance differentiates IBD patients and healthy individuals. Principal component analysis with health status as

sekwencjonowanie metagenomów

A map of diversity in the human microbiome

Lactobacillus species (*L. gasseri*, *L. jensenii*, *L. crispatus*, *L. iners*) are predominant but mutually exclusive in the **vagina**



Staphylococcus epidermidis colonizes external body sites



○ Commensal microbes
★ Potential pathogens

The four most abundant phyla

- Actinobacteria
- Bacteroidetes
- Firmicutes
- Proteobacteria

Low abundance phyla

- Chloroflexi
- Cyanobacteria
- Euryarchaeota
- Fusobacteria
- Lentisphaerae
- Spirochaetes
- Synergistetes
- Tenericutes
- Thermi
- Verrucomicrobia

National Institutes of Health
Human Microbiome Project

N. Segata & C. Huttenhower
http://huttenhower.sph.harvard.edu
Microbiome Online (Cantor and Holt 2012) from EuroMicrobiome

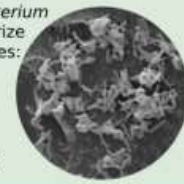


Streptococcus dominates the oral cavity with *S. mitis* > 75% in the **cheek**

Propionibacterium acnes lives on the skin and nose of most people



Many *Corynebacterium* species characterize different body sites:
C. matruchoti the **plaque**
C. accolens the **nose**
C. croppenstedtii the **skin**



Several *Prevotella* species are present in the gastrointestinal tract. *P. copri* is present in 19% of the subjects and dominates the **intestinal** flora when present



Microscopy from <http://barmap.wishartlab.com>

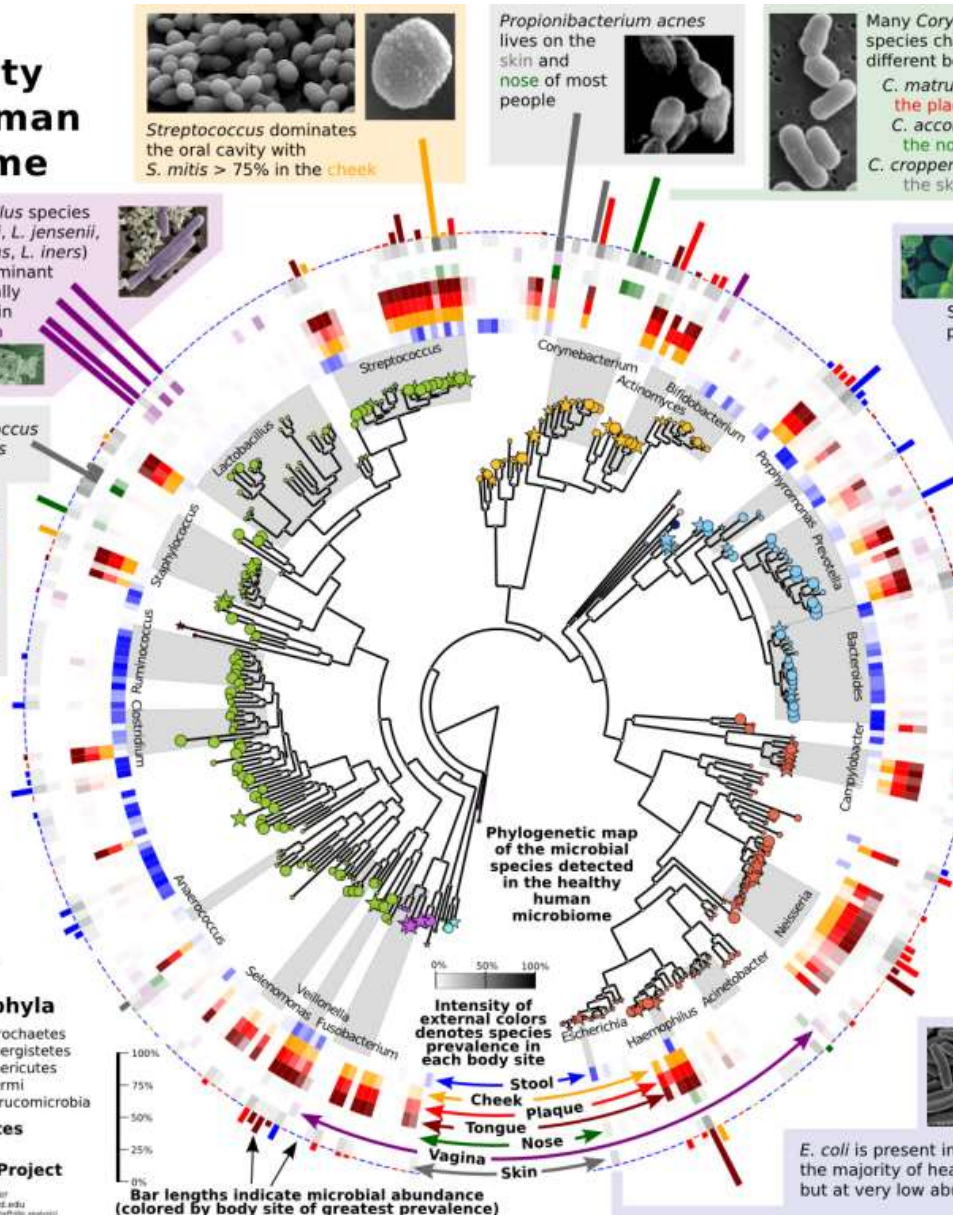
Bacteroides is the most abundant genus in the **gut** of almost all healthy subjects



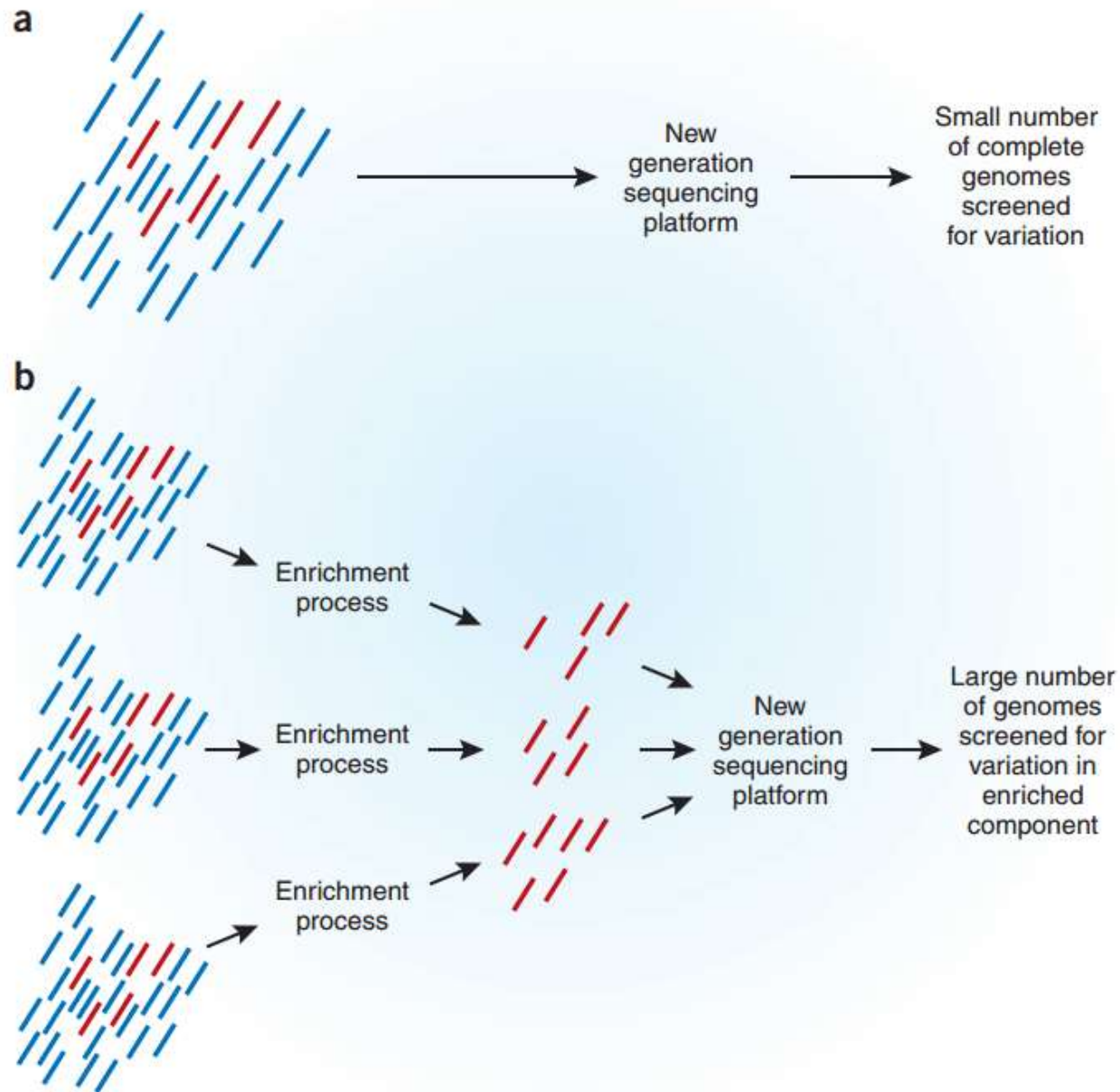
Campylobacter includes opportunistic pathogens, but members live in the oral cavities of most healthy people in the cohort



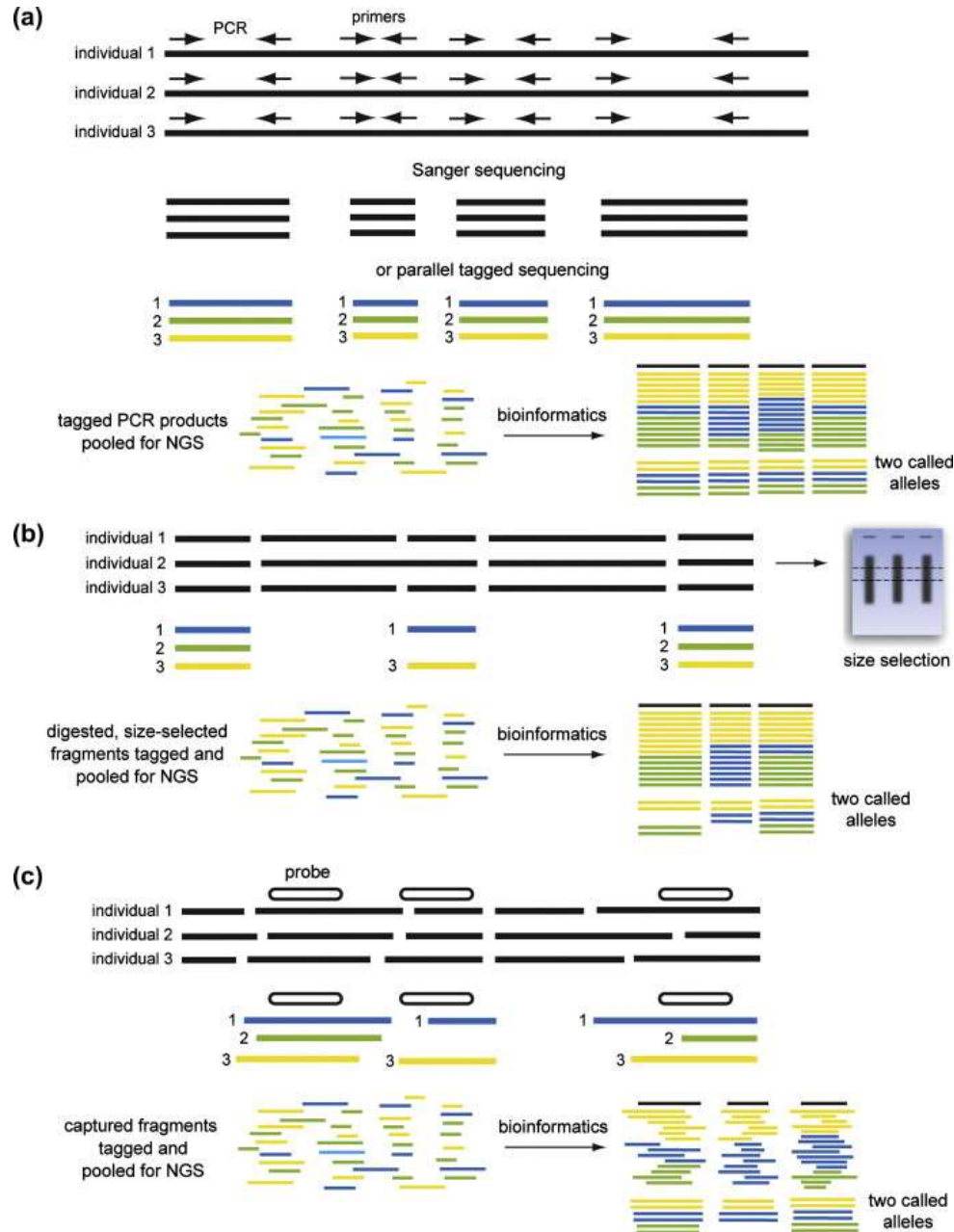
E. coli is present in the **gut** of the majority of healthy subjects but at very low abundance



resekwencjonowanie



resekwencjonowanie



resekwencjonowanie

Zalety:

- pozwala na uzyskanie tych samych informacji co sekwencjonowanie *de novo*
- wymaga znacznie mniejszego pokrycia sekwencji odczytami
- tańsze i szybsze

Wady:

- Składanie sekwencji wymaga genomu referencyjnego
- nie nadaje się do szybko zmieniających się genomów oraz organizmów o silnym zróżnicowaniu genetycznym w populacji

resekwencjonowanie

Resekwencjonowanie fragmentów genomów:

- pozwala na studiowanie istotnych fragmentów genomu w dużej skali (setki/tysiące osobników)
- najczęściej stosowane w identyfikacji zmienności genetycznej zasocjowanej z chorobami (rak, schorzenia genetyczne)

Rodzaje bibliotek zawierających fragmenty genomu:

- biblioteki egzomowe
- nakierowane na rejony uwikłane w schorzenia genetyczne
- nakierowane na uzupełnienie przerw w sekwencjach genomowych
- nakierowane na geny istotne dla analiz filogenetycznych

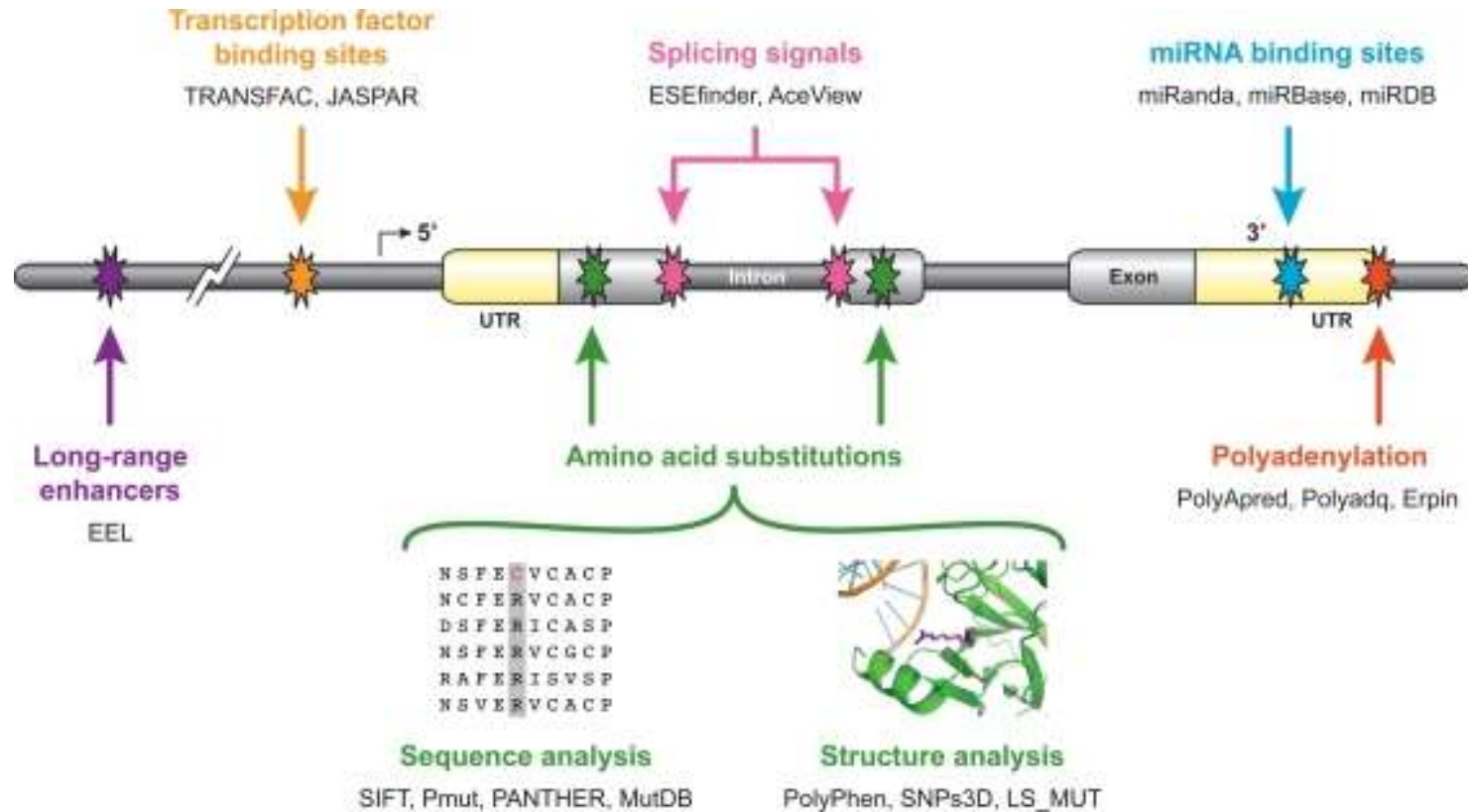
analiza wariacji

- SNP, krótkie delecje i insercje (pojedyncze odczyty, paired-ends)
- rearanżacje chromosomowe (odczyty paired-ends, mate-pair)
- wariacje liczby kopii (pojedyncze odczyty, paired-ends)

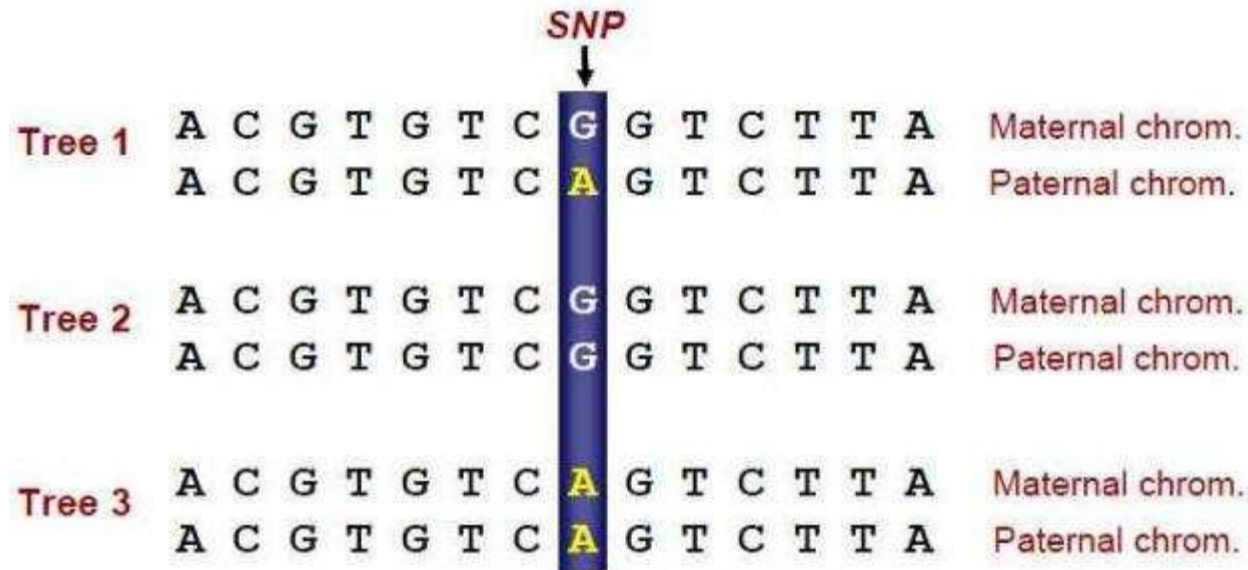
Podstawy:

- allele – wersje genu występująca w populacji (A, a)
- genotyp – zestawy alleli danego genu występujące w populacji (AA, Aa, aa)

analiza varijacije SNV/SNP



analiza varijacije SNV/SNP



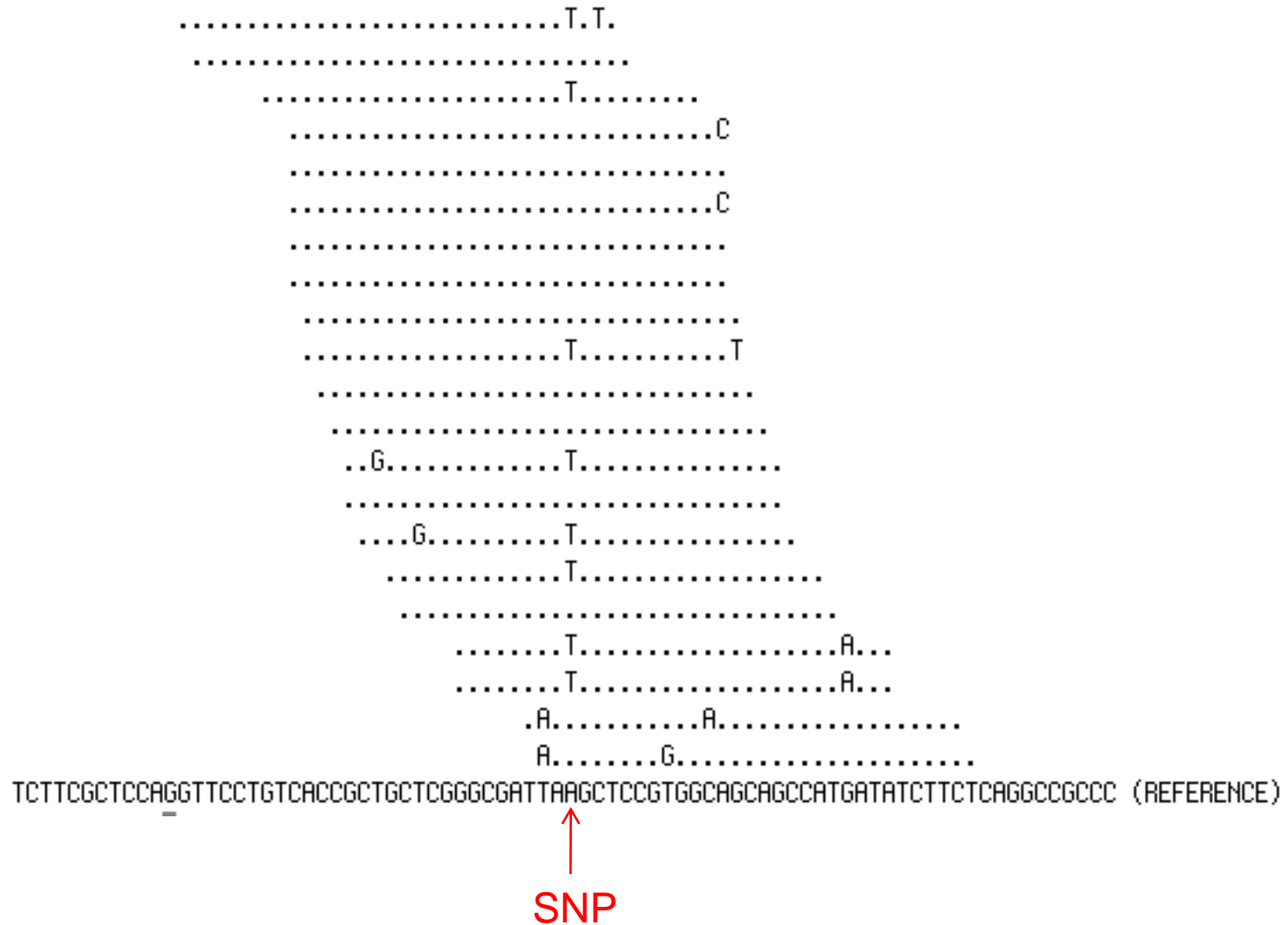
analiza wariacji SNV/SNP

```
C.....G...
C.....G...
.....G....
G.....G.....
.....GG.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....
.....G.....T.
A....G.....
..G.....
..G.....G.
G.....G...
```

AAATCCTGTAATTCAGGGTGATGCTGGTTTGACTGGACGCAAAATCATTGTGGACACTTATGGCGGTTGGGGTGCTCAT (REFERENCE)

↑
SNP

analiza wariacji SNV/SNP

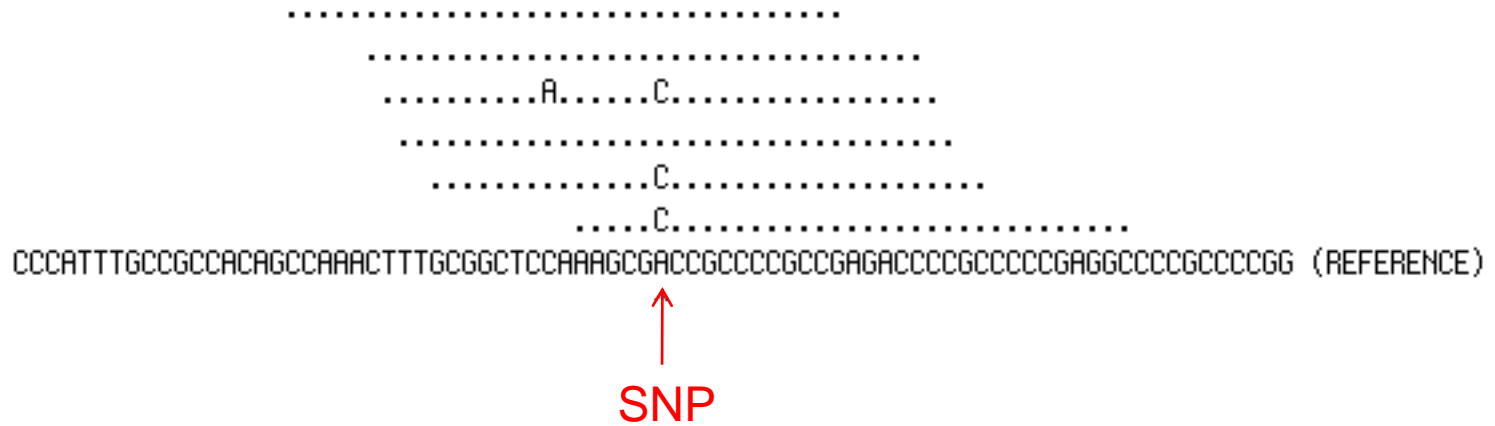


analiza wariacji SNV/SNP

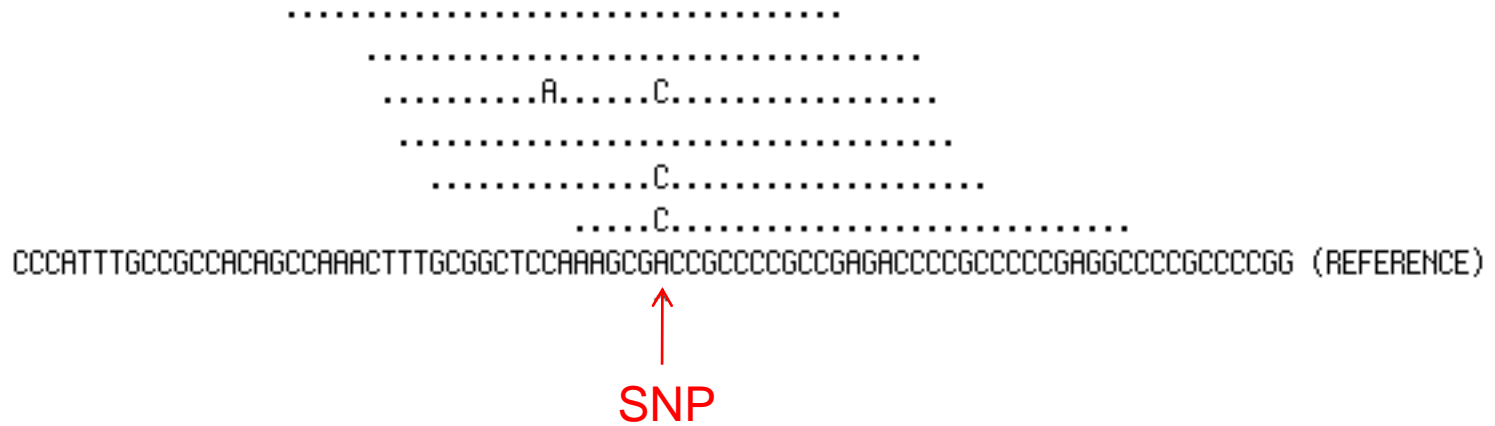
```
.....  
.....  
.....A.....C.....  
.....  
.....C.....  
.....C.....  
CCCATTTGCCGCCACGCCAACTTTGCGGCTCCAAGCGACCGCCCCGCCGAGACCCCGCCCCCGAGGCCCGCCCCGG (REFERENCE)
```

↑
SNP?

analiza wariacji SNV/SNP

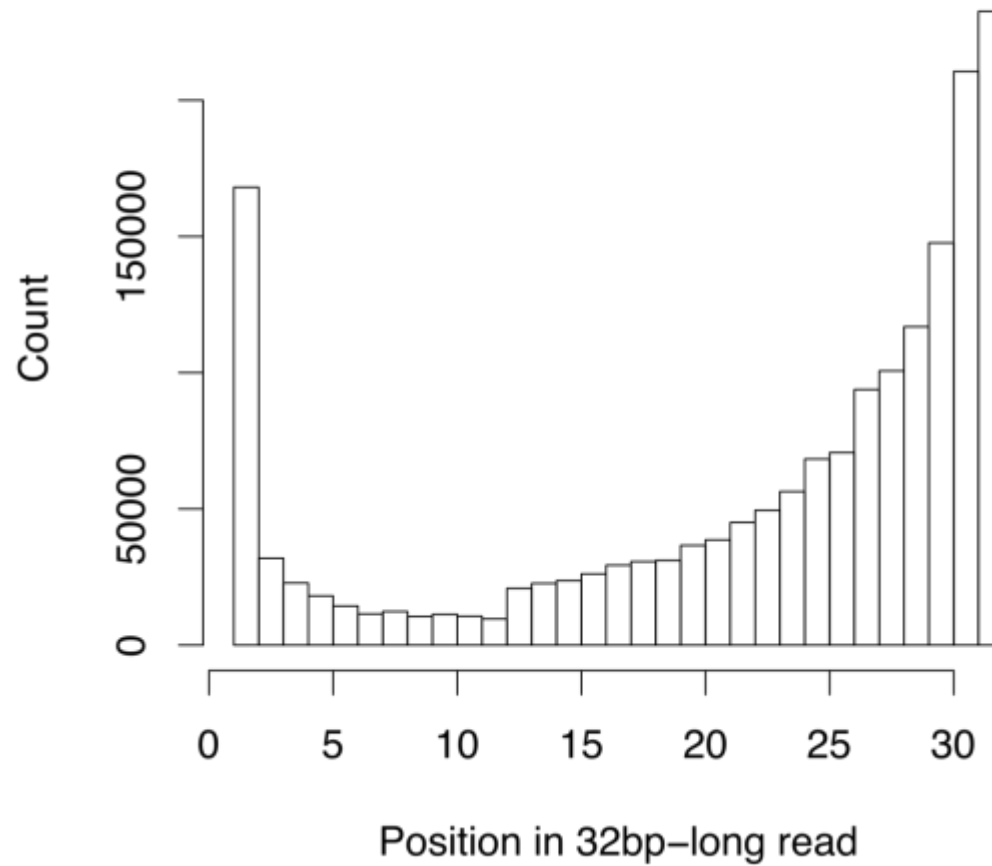


analiza wariacji SNV/SNP



analiza wariacji SNV/SNP

liczba błędów w sekwencji



genotypowanie

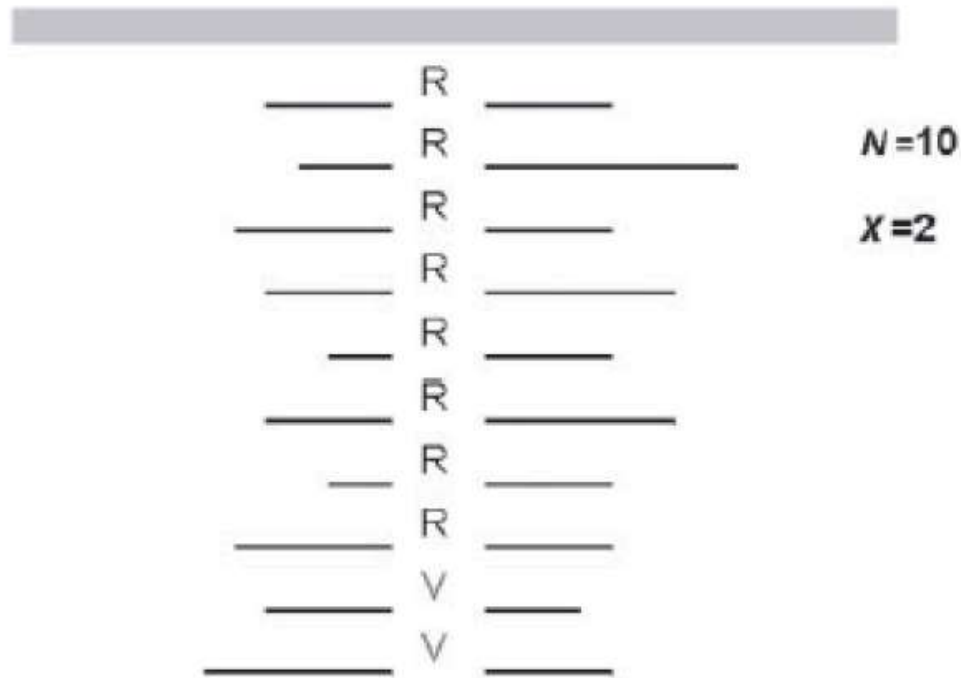


Fig. 1. Schematic of 10 aligned next-generation sequencing reads (R=reference nucleotide, V=variant nucleotide) for a single base position. N is read depth. X is variant count.

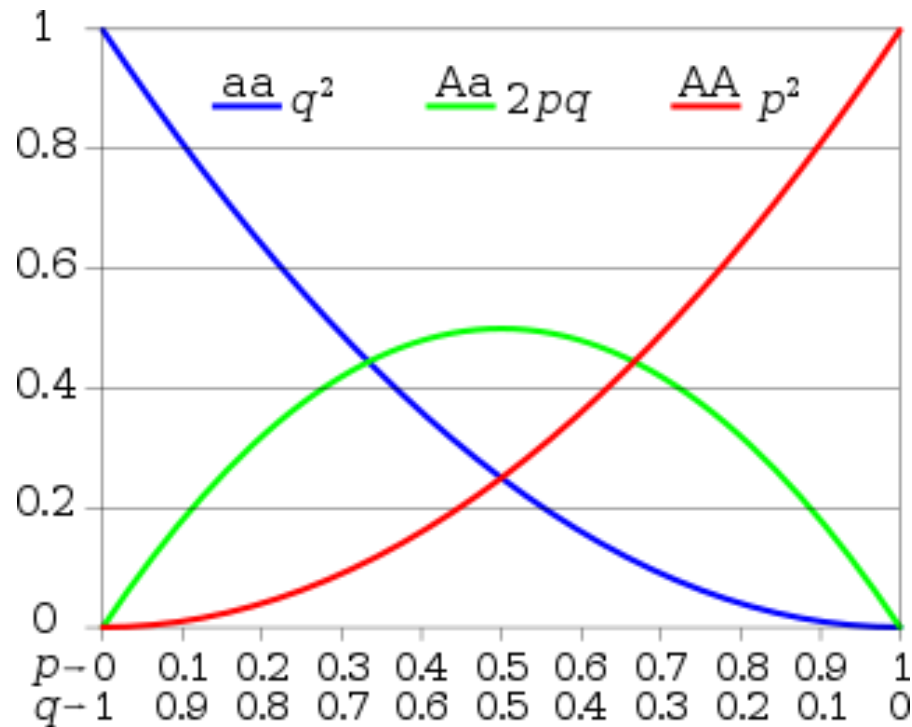
Cel: rozróżnić genotypy RR, RV i VV

- filtrowanie według jakości nukleotydów a następnie według zadanych progów:
 - jeśli % wariacji jest mniejszy niż próg 1 : RR
 - jeśli % wariacji jest większy niż próg 1 a mniejszy niż próg 2 : RV
 - jeśli % wariacji jest większy niż próg 2 : VV
- klasyfikator Bayesa
- Maksymalizacja oczekiwanych wariacji (expectation-maximalization) – optymalizacja parametrów aby maksymalizowały prawdopodobieństwo wariacji w obserwowanych danych

„state of the art” : Bayes

klasyfikator bayesa

- zakładamy równowagę Hardy'ego-Weinberg'a:
 - częstości alleli i genotypów w populacji są stałe w przypadku braku presji ewolucyjnej



klasyfikator bayesa

- zakładamy równowagę Hardy'ego-Weinberg'a:
 - częstości alleli i genotypów w populacji są stałe w przypadku braku presji ewolucyjnej
- Punkt wyjściowy: oczekiwane częstości i częstości błędów sekwencjonowania
- obliczamy wstępne prawdopodobieństwo dla genotypów RR, RV, VV
- szukamy :

$$P(G|D) = \frac{P(D|G)P(G)}{P(D)} = \frac{P(D|G) P(G)}{\sum_{i=1}^n P(D|G_i) P(G_i)}$$

częstości:

G = wstępne dla testowanego genotypu

D = obserwowane dla danych

- używamy $P(G|D)$ jako końcowego prawdopodobieństwa dla $P(G)$

- korzysta z algorytmu MapReduce (Google, Yahoo!)
- nieograniczona paralelizacja, gdyż identyfikacja SNP i genotypowanie odbywa się dla każdego locus niezależnie
- może analizować oraz łączyć wiele próbek
- używa uniwersalnego formatu bam

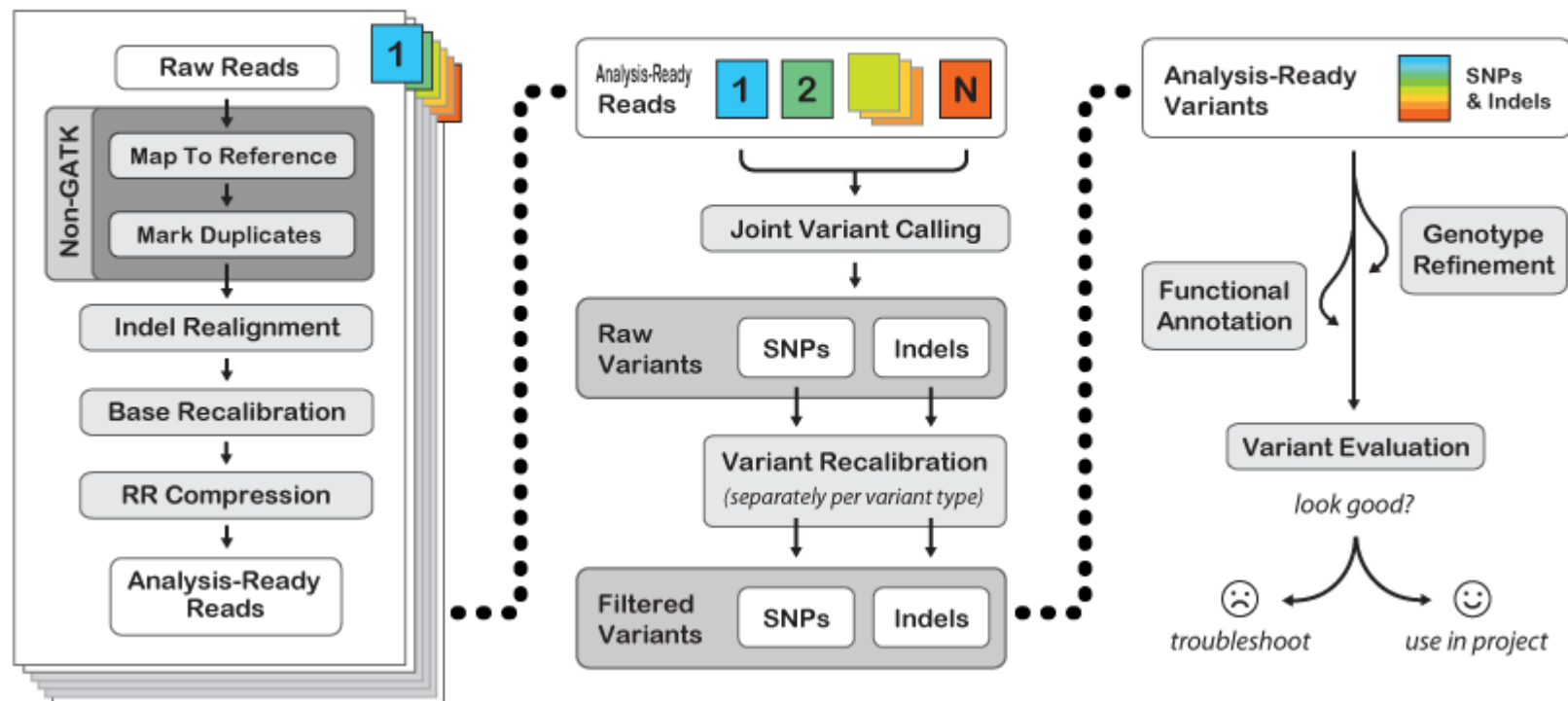
Data Pre-processing

>>

Variant Discovery

>>

Preliminary Analyses



realignment

wiele porównań do referencji ≠ wielokrotne porównanie

```
ref  aggtttataaaac----aattaagtctacagagcaacta
sample aggtttataaaacAAATaattaagtctacagagcaacta
read1 aggtttataaaac****aaAtaa
read2  ggttttataaaac****aaAtaaTt
read3   ttataaaacAAATaattaagtctaca
read4    CaaaT****aattaagtctacagagcaac
read5     aaT****aattaagtctacagagcaact
read6      T****aattaagtctacagagcaacta
```

```
ref  aggtttataaaac----aattaagtctacagagcaacta
sample aggtttataaaacAAATaattaagtctacagagcaacta
read1 aggtttataaaacAAATaa
read2  ggttttataaaacAAATaatt
read3   ttataaaacAAATaattaagtctaca
read4    cAAATaattaagtctacagagcaac
read5     AATaattaagtctacagagcaact
read6      Taattaagtctacagagcaacta
```

realignment

wiele porównań do referencji ≠ wielokrotne porównanie

```
ref   aggtttataaaacAAAAaattaagtctacagagcaacta
sample aggtttataaaacAAA-aattaagtctacagagcaacta
read1  aggtttataaaacAA-Aaattaagtctacagagcaacta
read2  aggtttataaaacA-AAaattaagtctacagagcaacta
read3  aggtttataaaac-AAaattaagtctacagagcaacta
consensus aggtttataaaacAAAAaattaagtctacagagcaacta
```

Rozwiązania:

- realignment rejonów zawierających liczne wyspy/SNP (GATK)
- wyznaczenie prawdopodobieństwa że pozycja jest źle dopasowana i usunięcie problematycznych rejonów z analizy (SAMTOOLS)

rekalibracja wariacji

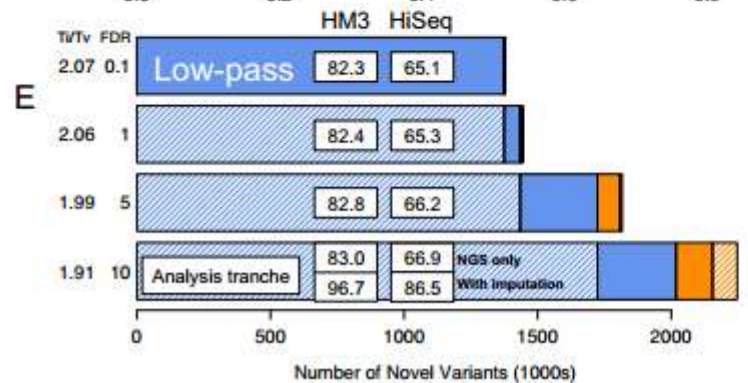
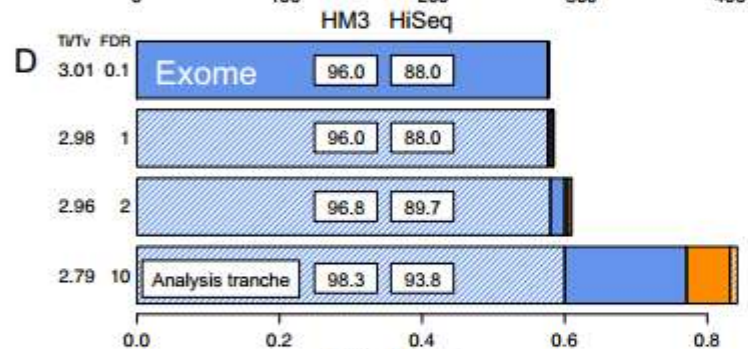
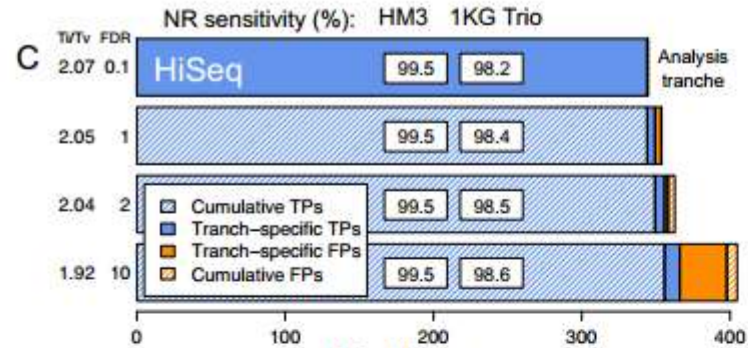
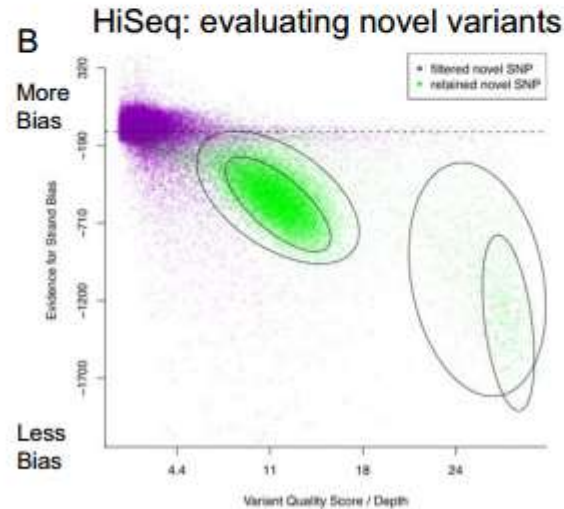
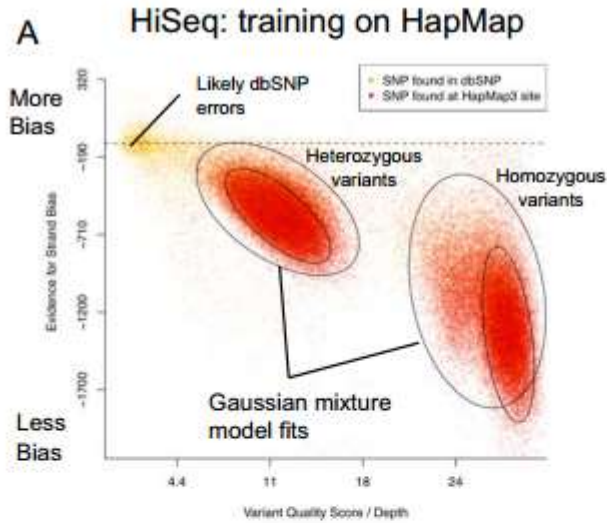
Dlaczego:

- programy zwracają dużo fałszywych wariacji
- liczba faktycznych wariacji genetycznych zależy od jakości i ilości konkretnych danych

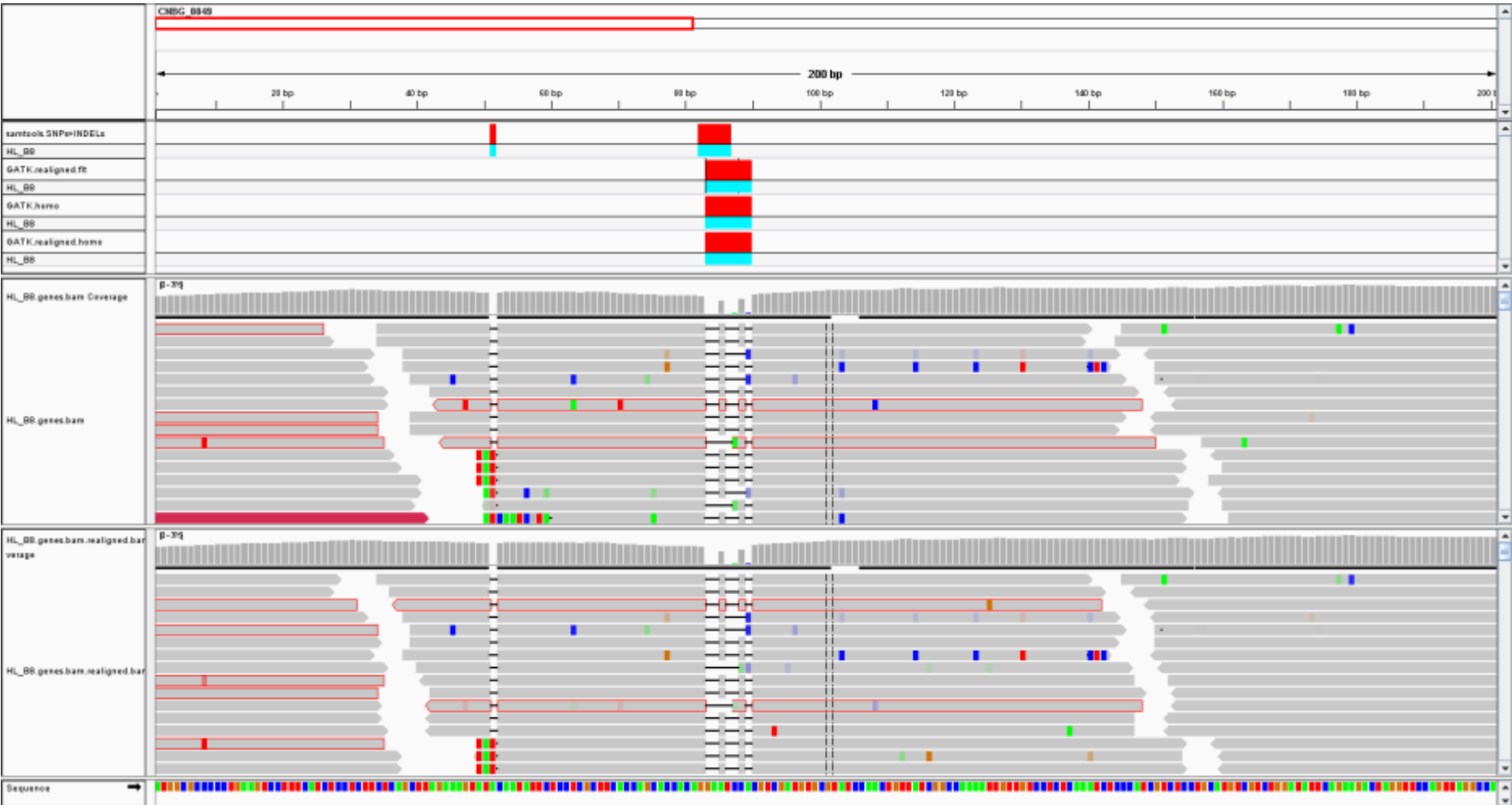
Jak:

- znane wariacje służą jako wzorzec do wyznaczenia prawdopodobieństwa dla nowych wariacji (np. z HapMap)

rekalibracija varijacije



analiza wariacji krótkie insercje/delecje (indels)

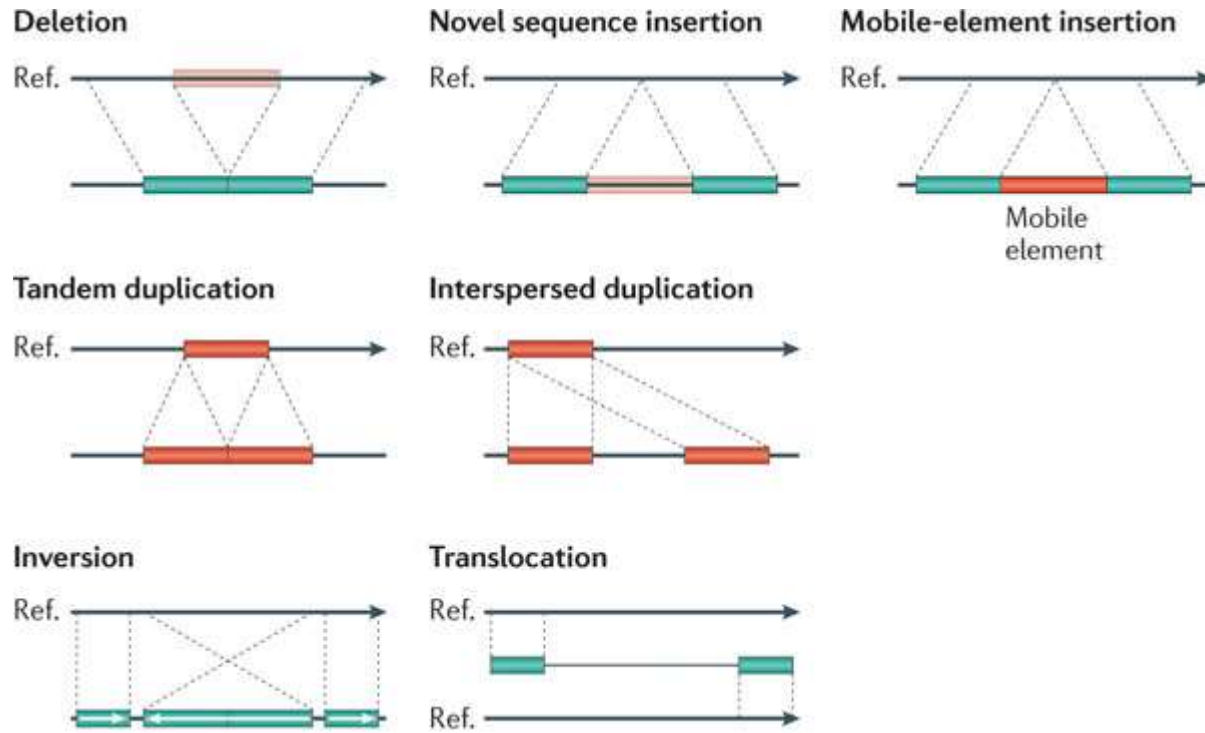


analiza wariacji SNP/indels

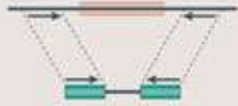


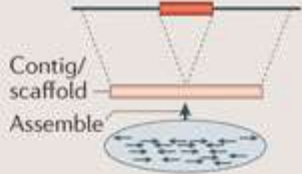
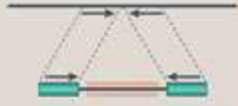
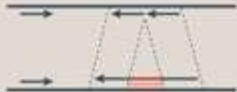
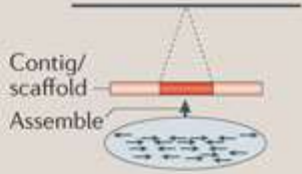
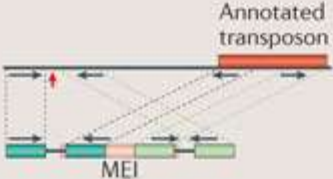
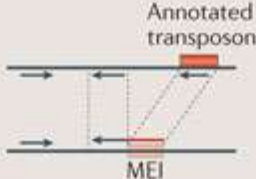
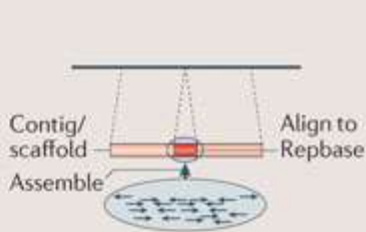
Narzędzia:

1. GATK Unified Genotyper
2. Komenda MPILEUPz pakietu SAMTOOLS

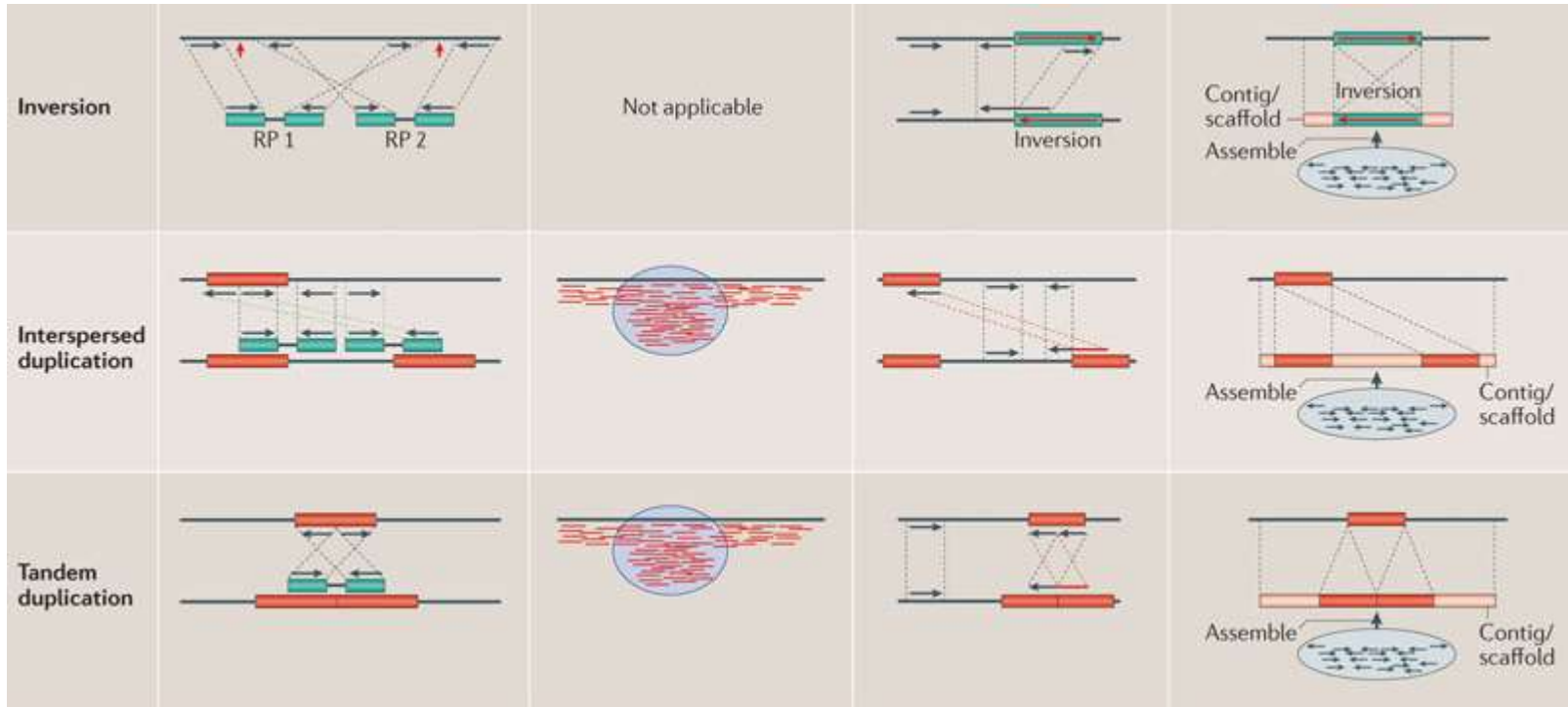
analiza varijacije rearanžacije chromosomalne



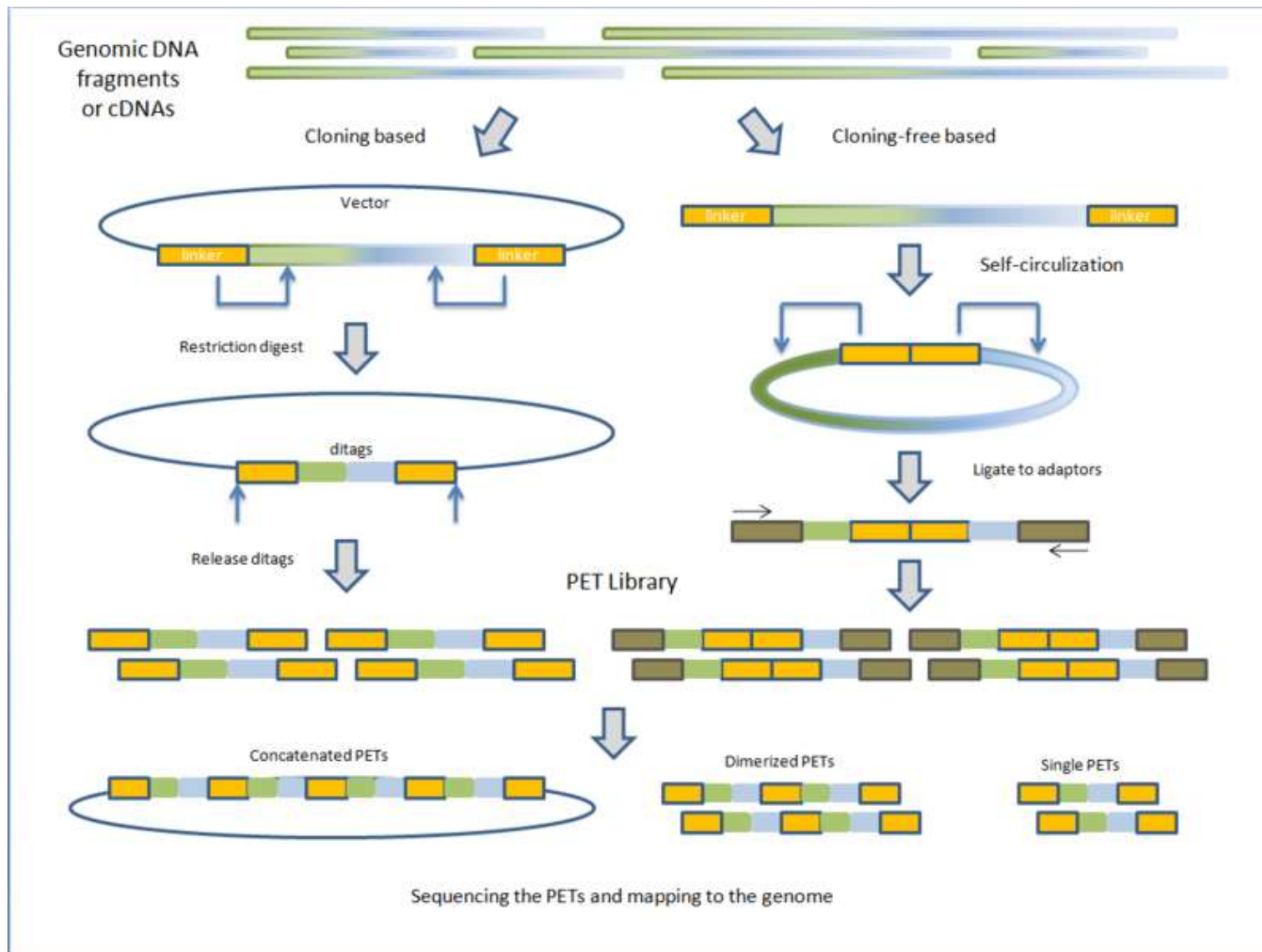
analiza varijacije rearanžacije chromosomalne

SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Mobile-element insertion		Not applicable		

analiza varijacije rearanžacije chromosomalne

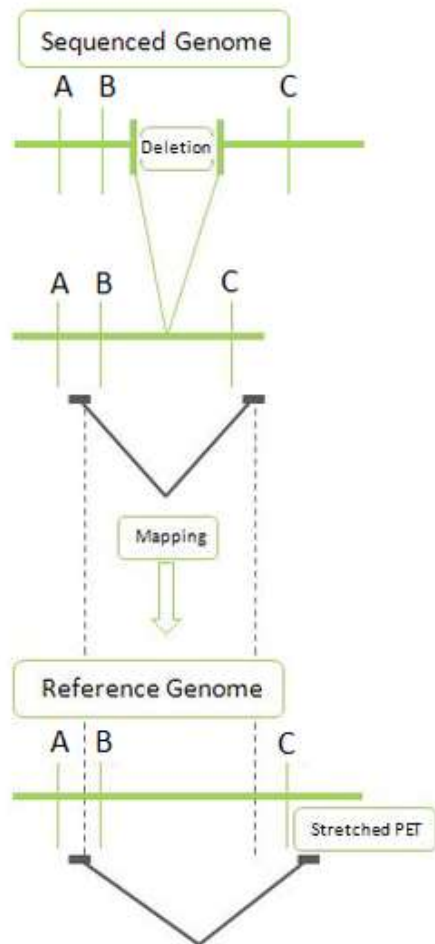


analiza wariacji rearanżacje chromosomalne – PET sequencing

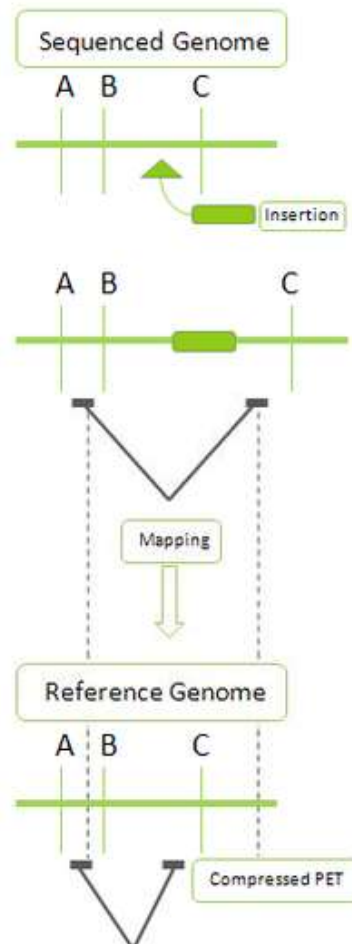


analiza wariacji rearanżacje chromosomalne – PET sequencing

Deletion Detection with PET



Insertion Detection with PET



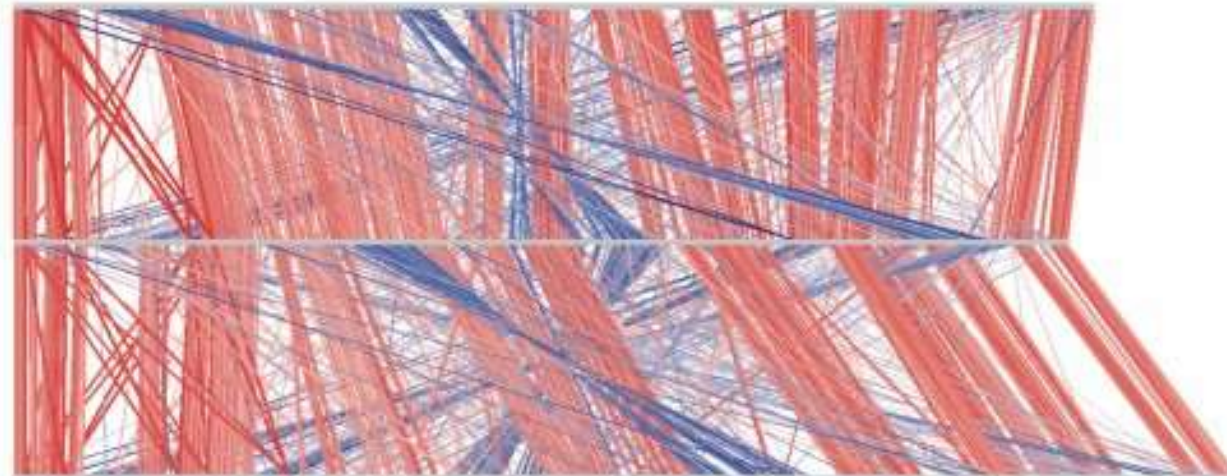
analiza varijacije rearanžacije chromosomalne

A

S. pyogenes MGAS315

S. uberis 0140J

S. zooepidemicus H70



0 Mb

1.0 Mb

2.0 Mb

B

S. uberis 0140J

S. agalactiae NEM31E



0 Mb

1.0 Mb

VisualLightBox.com

The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups

Christina Curtis^{1,2†*}, Sohrab P. Shah^{3,4*}, Suet-Feung Chin^{1,2*}, Gulisa Turashvili^{3,4*}, Oscar M. Rueda^{1,2}, Mark J. Dunning², Doug Speed^{2,5†}, Andy G. Lynch^{1,2}, Shamith Samarajiwa^{1,2}, Yinyin Yuan^{1,2}, Stefan Gräßl^{1,2}, Gavin Ha³, Gholamreza Haffari³, Ali Bashashati³, Roslin Russell², Steven McKinney^{3,4}, METABRIC Group[†], Anita Langerød⁶, Andrew Green⁷, Elena Provenzano⁸, Gordon Wishart⁸, Sarah Pinder⁹, Peter Watson^{3,4,10}, Florian Markowitz^{1,2}, Leigh Murphy¹⁰, Ian Ellis⁷, Arnie Purushotham^{9,11}, Anne-Lise Børresen-Dale^{6,12}, James D. Brenton^{2,13}, Simon Tavaré^{1,2,5,14}, Carlos Caldas^{1,2,8,13} & Samuel Aparicio^{3,4}

The elucidation of breast cancer subgroups and their molecular drivers requires integrated views of the genome and transcriptome from representative numbers of patients. We present an integrated analysis of copy number and gene expression in a discovery and validation set of 997 and 995 primary breast tumours, respectively, with long-term clinical follow-up. Inherited variants (copy number variants and single nucleotide polymorphisms) and acquired somatic copy number aberrations (CNAs) were associated with expression in ~40% of genes, with the landscape dominated by *cis*- and *trans*-acting CNAs. By delineating expression outlier genes driven in *cis* by CNAs, we identified putative cancer genes, including deletions in *PPP2R2A*, *MTAP* and *MAP2K4*. Unsupervised analysis of paired DNA–RNA profiles revealed novel subgroups with distinct clinical outcomes, which reproduced in the validation cohort. These include a high-risk, oestrogen-receptor-positive 11q13/14 *cis*-acting subgroup and a favourable prognosis subgroup devoid of CNAs. *Trans*-acting aberration hotspots were found to modulate subgroup-specific gene networks, including a TCR deletion-mediated adaptive immune response in the ‘CNA-devoid’ subgroup and a basal-specific chromosome 5 deletion-associated mitotic network. Our results provide a novel molecular stratification of the breast cancer population, derived from the impact of somatic CNAs on the transcriptome.