

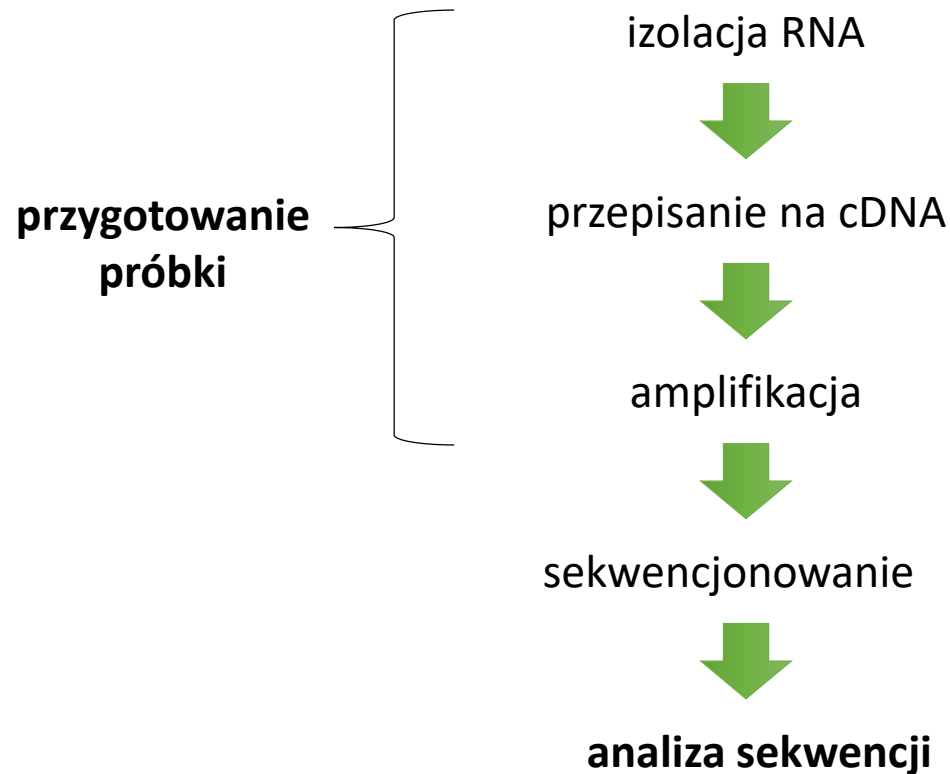
A hand is holding a clear glass test tube against a blue background. Inside the test tube, a DNA double helix is visible, with the two strands represented by silver ribbons and the base pairs by colorful rectangular blocks in shades of red, blue, purple, and orange. The DNA structure is positioned diagonally within the tube.

Analiza danych z wysokoprzepustowego  
sekwencjonowania

## **W2: Sekwencjonowanie transkryptomu**

# sekwencjonowanie transkryptomu

## Transkryptomika



# sekwencjonowanie transkryptomu

## Zalety:

- metoda bezpośrednia – widzimy to co jest, a nie to o czym wiemy że jest
- metoda skalowalna – możemy dopasować do swoich potrzeb
- metoda adaptowalna – możliwe badanie różnych zjawisk poprzez odpowiednie przygotowanie biblioteki cDNA

## Wady:

- wyższy koszt niż w przypadku mikromacierzy

# profilowanie ekspresji genów

## A Global View of Gene Activity and Alternative Splicing by Deep Sequencing of the Human Transcriptome

Marc Sultan,<sup>1\*</sup> Marcel H. Schulz,<sup>2,3\*</sup> Hugues Richard,<sup>2\*</sup> Alon Magen,<sup>1</sup> Andreas Klingenhoff,<sup>4</sup> Matthias Scherf,<sup>4</sup> Martin Seifert,<sup>4</sup> Tatjana Borodina,<sup>1</sup> Aleksey Soldatov,<sup>1</sup> Dmitri Parkhomchuk,<sup>1</sup> Dominic Schmidt,<sup>1</sup> Sean O'Keefe,<sup>2</sup> Stefan Haas,<sup>2</sup> Martin Vingron,<sup>2</sup> Hans Lehrach,<sup>1</sup> Marie-Laure Yaspo<sup>1†</sup>

The functional complexity of the human transcriptome is not yet fully elucidated. We report a high-throughput sequence of the human transcriptome from a human embryonic kidney and a B cell line. We used shotgun sequencing of transcripts to generate randomly distributed reads. Of these, 50% mapped to unique genomic locations, of which 80% corresponded to known exons. We found that 66% of the polyadenylated transcriptome mapped to known genes and 34% to nonannotated genomic regions. On the basis of known transcripts, RNA-Seq can detect 25% more genes than can microarrays. A global survey of messenger RNA splicing events identified 94,241 splice junctions (4096 of which were previously unidentified) and showed that exon skipping is the most prevalent form of alternative splicing.

## Comparative Functional Genomics of the Fission Yeasts

Nicholas Rhind,<sup>1¶</sup> Zehua Chen,<sup>2</sup> Moran Yassour,<sup>3,4,5¶</sup> Dawn A. Thompson,<sup>3¶</sup> Brian J. Haas,<sup>2¶</sup> Naomi Habib,<sup>5,6¶</sup> Ilan Wapinski,<sup>3,7¶</sup> Sushmita Roy,<sup>3,8¶</sup> Michael F. Lin,<sup>8</sup> David I. Heiman,<sup>2</sup> Sarah K. Young,<sup>2</sup> Kanji Furuya,<sup>9</sup> Yabin Guo,<sup>10</sup> Alison Pidoux,<sup>11</sup> Huei Mei Chen,<sup>12</sup> Barbara Robbette,<sup>13\*</sup> Jonathan M. Goldberg,<sup>2</sup> Keita Aoki,<sup>9</sup> Elizabeth H. Bayne,<sup>11†</sup> Aaron M. Berlin,<sup>2</sup> Christopher A. Desjardins,<sup>2</sup> Edward Dobbs,<sup>11</sup> Livio Dukaj,<sup>1</sup> Lin Fan,<sup>2</sup> Michael G. FitzGerald,<sup>2</sup> Courtney French,<sup>6</sup> Sharvari Gujja,<sup>2</sup> Klavs Hansen,<sup>14†</sup> Dan Keifenheim,<sup>1</sup> Joshua Z. Levin,<sup>2</sup> Rebecca A. Mosher,<sup>15§</sup> Carolin A. Müller,<sup>16</sup> Jenna Pfiffner,<sup>2</sup> Margaret Priest,<sup>2</sup> Carsten Russ,<sup>2</sup> Agata Smialowska,<sup>17,18</sup> Peter Swoboda,<sup>17</sup> Sean M. Sykes,<sup>2</sup> Matthew Vaughn,<sup>14</sup> Sonya Vengrova,<sup>19</sup> Ryan Yoder,<sup>13</sup> Qiandong Zeng,<sup>2</sup> Robin Allshire,<sup>11</sup> David Baulcombe,<sup>15</sup> Bruce W. Birren,<sup>20</sup> William Brown,<sup>16</sup> Karl Ekwall,<sup>17,18</sup> Manolis Kellis,<sup>8,3</sup> Janet Leatherwood,<sup>12</sup> Henry Levin,<sup>10</sup> Hanah Margalit,<sup>6</sup> Rob Martienssen,<sup>14</sup> Conrad A. Nieduszynski,<sup>16</sup> Joseph W. Spatafora,<sup>13</sup> Nir Friedman,<sup>5,21</sup> Jacob Z. Dalgaard,<sup>19</sup> Peter Baumann,<sup>22,23,24</sup> Hironori Niki,<sup>9</sup> Aviv Regev,<sup>3,4,24¶</sup> Chad Nusbaum<sup>2¶</sup>

## The Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing

Ugrappa Nagalakshmi,<sup>1\*</sup> Zhong Wang,<sup>1\*</sup> Karl Waern,<sup>2</sup> Chong Shou,<sup>2</sup> Debasish Raha,<sup>1</sup> Mark Gerstein,<sup>2,3</sup> Michael Snyder<sup>1,2,3†</sup>

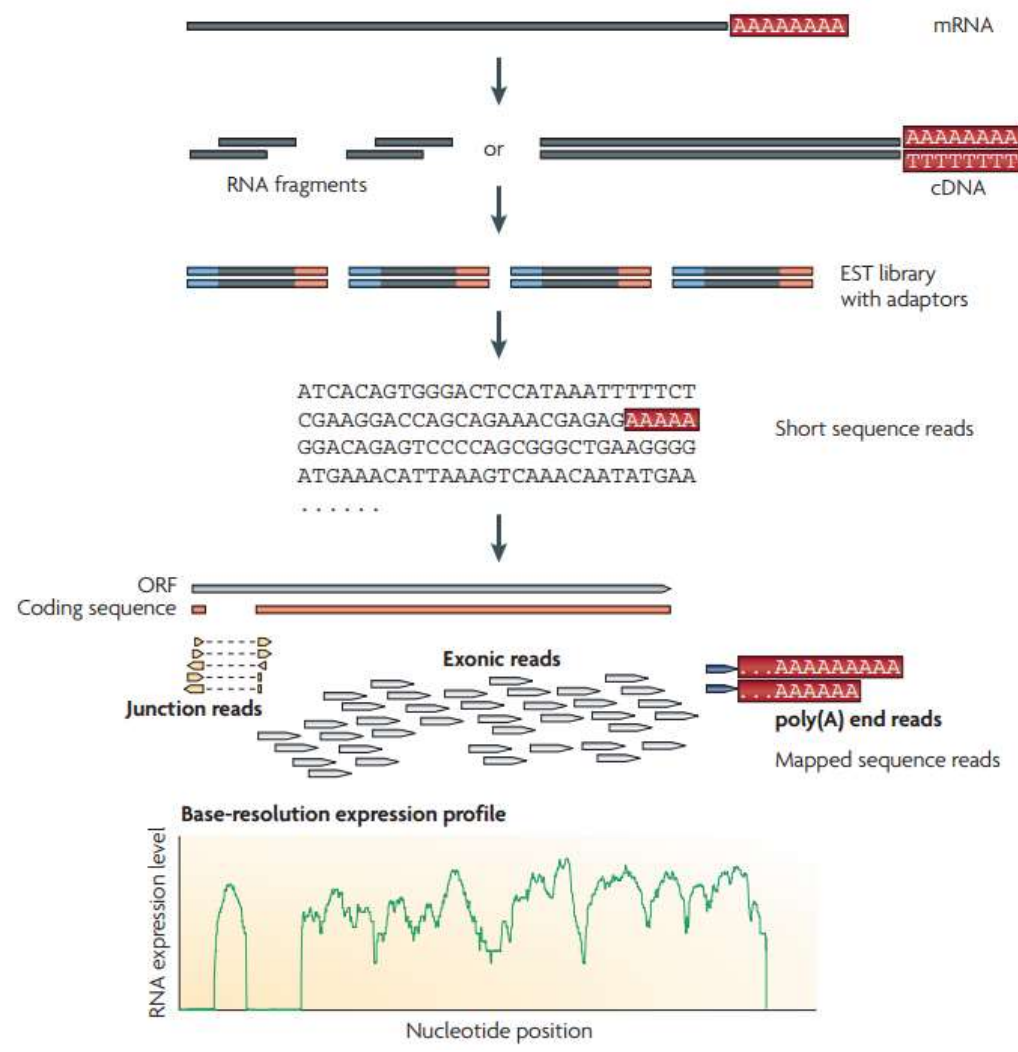
The identification of untranslated regions, introns, and coding regions within an organism remains challenging. We developed a quantitative sequencing-based method called RNA-Seq for mapping transcribed regions, in which complementary DNA fragments are subjected to high-throughput sequencing and mapped to the genome. We applied RNA-Seq to generate a high-resolution transcriptome map of the yeast genome and demonstrated that most (74.5%) of the nonrepetitive sequence of the yeast genome is transcribed. We confirmed many known and predicted introns and demonstrated that others are not actively used. Alternative initiation codons and upstream open reading frames also were identified for many yeast genes. We also found unexpected 3'-end heterogeneity and the presence of many overlapping genes. These results indicate that the yeast transcriptome is more complex than previously appreciated.

## Cold-Inducible RNA-Binding Protein Modulates Circadian Gene Expression Posttranscriptionally

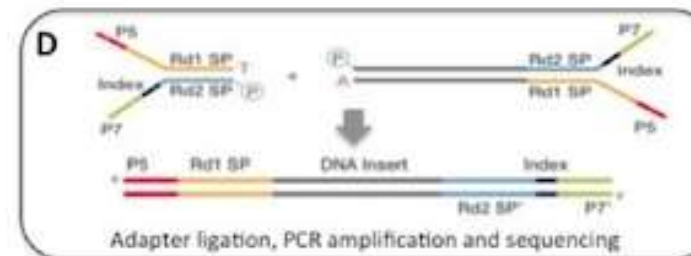
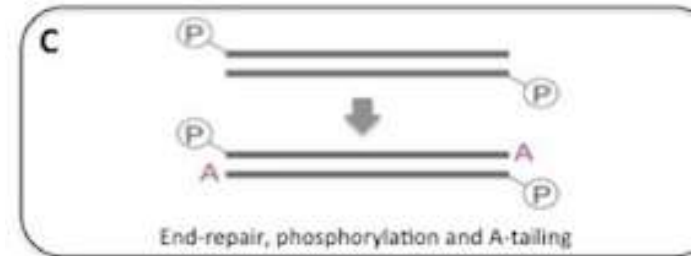
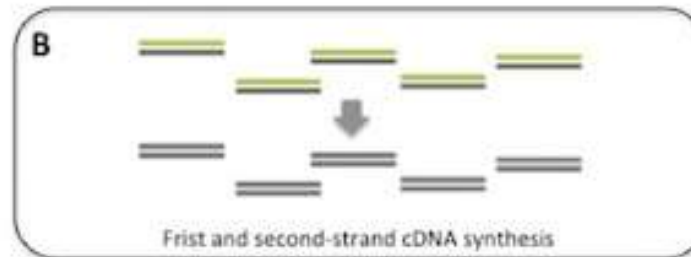
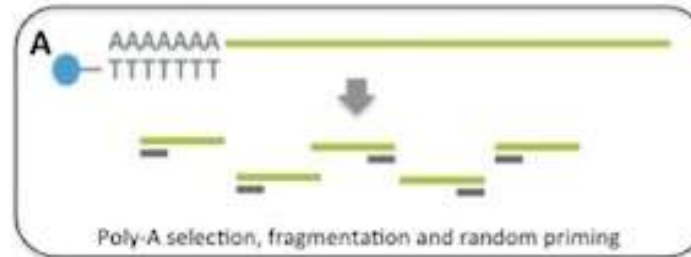
Jörg Morf,<sup>1</sup> Guillaume Rey,<sup>2\*</sup> Kim Schneider,<sup>2</sup> Markus Stratmann,<sup>1</sup> Jun Fujita,<sup>3</sup> Felix Naef,<sup>2</sup> Ueli Schibler<sup>1†</sup>

In mammalian tissues, circadian gene expression can be driven by local oscillators or systemic signals controlled by the master pacemaker in the suprachiasmatic nucleus. We show that simulated body temperature cycles, but not peripheral oscillators, controlled the rhythmic expression of cold-inducible RNA-binding protein (CIRP) in cultured fibroblasts. In turn, loss-of-function experiments indicated that CIRP was required for high-amplitude circadian gene expression. The transcriptome-wide identification of CIRP-bound RNAs by a biotin-streptavidin-based cross-linking and immunoprecipitation (CLIP) procedure revealed several transcripts encoding circadian oscillator proteins, including CLOCK. Moreover, CLOCK accumulation was strongly reduced in CIRP-depleted fibroblasts. Because ectopic expression of CLOCK improved circadian gene expression in these cells, we surmise that CIRP confers robustness to circadian oscillators through regulation of CLOCK expression.

# profilowanie ekspresji genów - RNA-seq



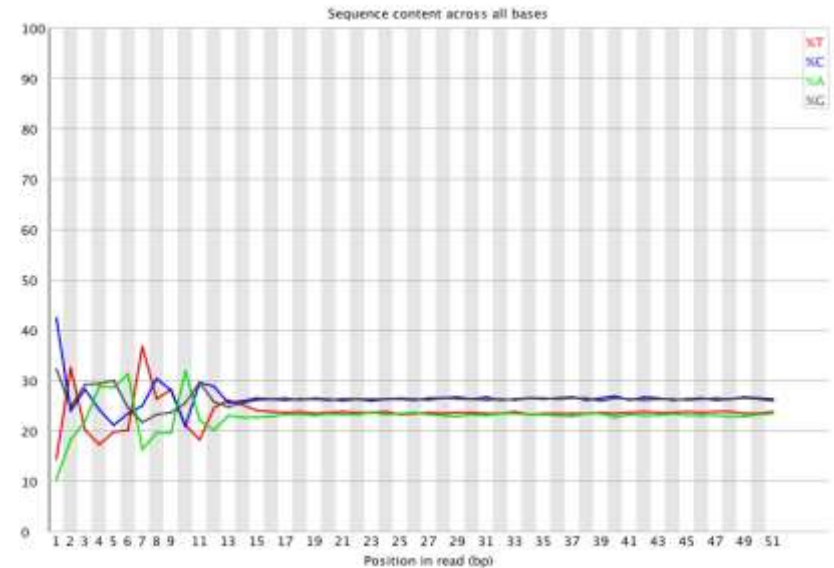
# profilowanie ekspresji genów - RNA-seq illumina Truseq



# profilowanie ekspresji genów - RNA-seq illumina Truseq

## Znaczenie dla analizy:

- możliwość występowania preferencji w rozkładzie nukleotydów na początku sekwencji (związane z losowymi starterami)
- metoda daje dobrą „losowość” fragmentacji - występowanie duplikatów spowodowane jest z dużym prawdopodobieństwem reakcją PCR
- fragmentacja i selekcja odpowiedniej długości insertów powoduje prawie całkowity brak 3' adaptorów
- utrata informacji o nici z której pochodzi transkrypt – podczas zliczania odczytów należy uwzględnić obie nici



# profilowanie ekspresji genów - RNA-seq illumina Truseq

## Zalety:

- Minimalna ilość manipulacji próbką przed amplifikacją
- Wysoka dokładność w oznaczaniu poziomu ekspresji genów
- Wysoka jakość biblioteki
- Wysoka dokładność w identyfikacji SNP

## Wady:

- Utrata informacji o nici z której pochodzą transkrypty

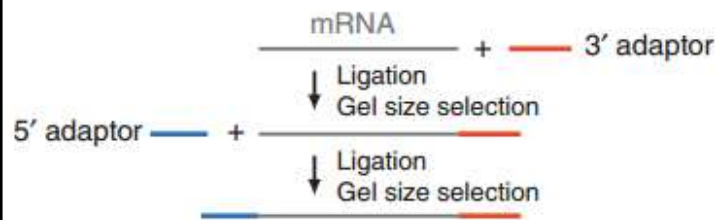
## Zastosowanie:

- Badanie poziomu ekspresji GENÓW
- Badanie aktywności transkrypcyjnej genomu
- Badanie polimorfizmu nukleotydowego na poziomie transkryptomu

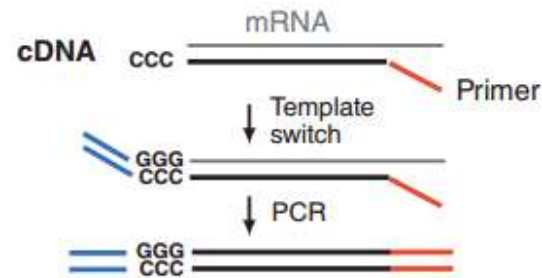


# profilowanie ekspresji genów - RNA-seq wprowadzenie adaptorów

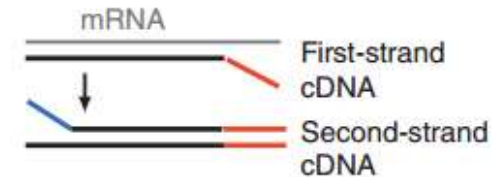
ligacja adaptorów do RNA



SMART



losowe primery z adaptorami



# profilowanie ekspresji genów - RNA-seq wprowadzenie adaptorów

## Znaczenie dla analizy:

- Utrzymanie informacji o kierunku transkrypcji – przy zliczaniu można rozróżnić geny kodowane na przeciwnych niciach
- Trakty homonukleotydowe dodawane w metodzie SMART bywają trudne do rozróżnienia od naturalnych traktów homopolimerowych
- Możliwe do zaobserwowania preferencje sekwencyjne na obu końcach odczytów związane z preferencjami ligazy

# profilowanie ekspresji genów - RNA-seq wprowadzenie adaptorów

## Zalety:

- Utrzymanie informacji o kierunku transkrypcji
- Możliwe uzyskanie wariantów splicingowych w skomplikowanych genomach

## Wady:

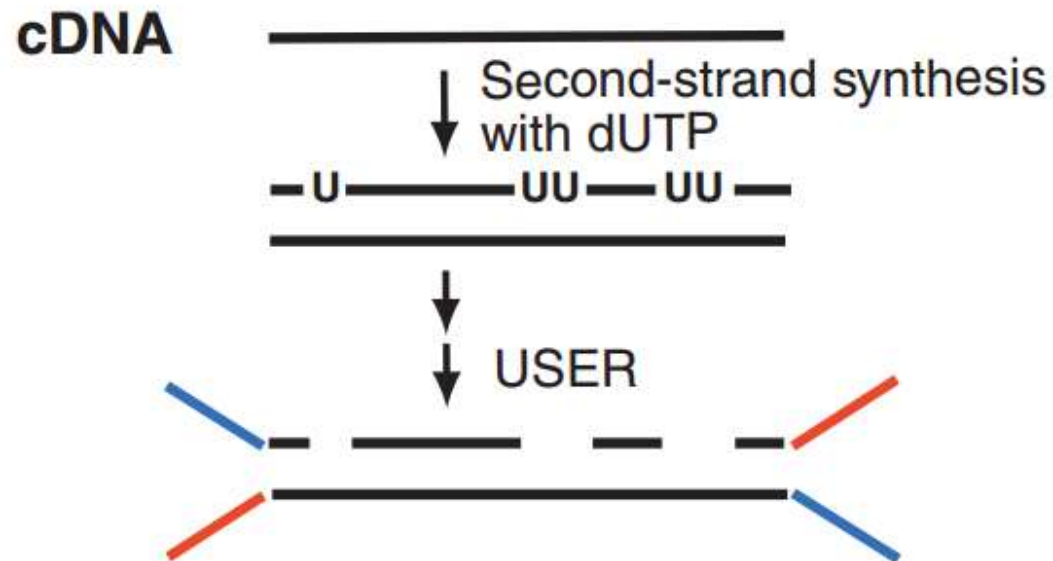
- Zaburzenia odczytywanego poziomu transkryptów poprzez użycie adaptorów
- Stosunkowo niska reproduktywność wyników
- Utrudniona analiza sekwencji

## Zastosowanie:

- Badanie poziomu ekspresji GENÓW i TRANSKRYPTÓW
- Badanie aktywności transkrypcyjnej genomu
- Badanie struktury genów
- Identyfikacja izoform splicingowych

# profilowanie ekspresji genów - RNA-seq dUTP

usuwanie drugiej nici cDNA z zastosowaniem dUTP



# profilowanie ekspresji genów - RNA-seq dUTP

## Zalety:

- Utrzymanie informacji o kierunku transkrypcji
- Możliwe uzyskanie wariantów splicingowych w skomplikowanych genomach
- Jakość biblioteki i powtarzalność wyników porównywalna z Truseq

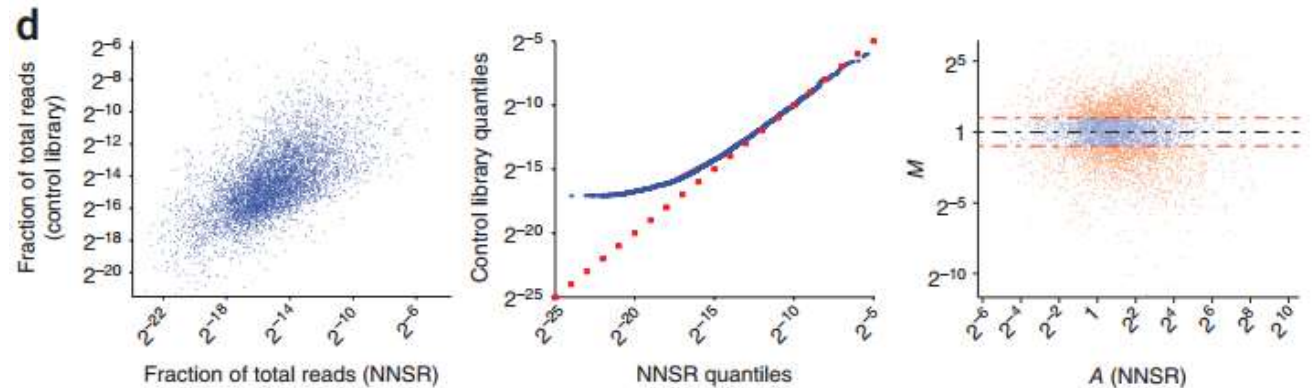
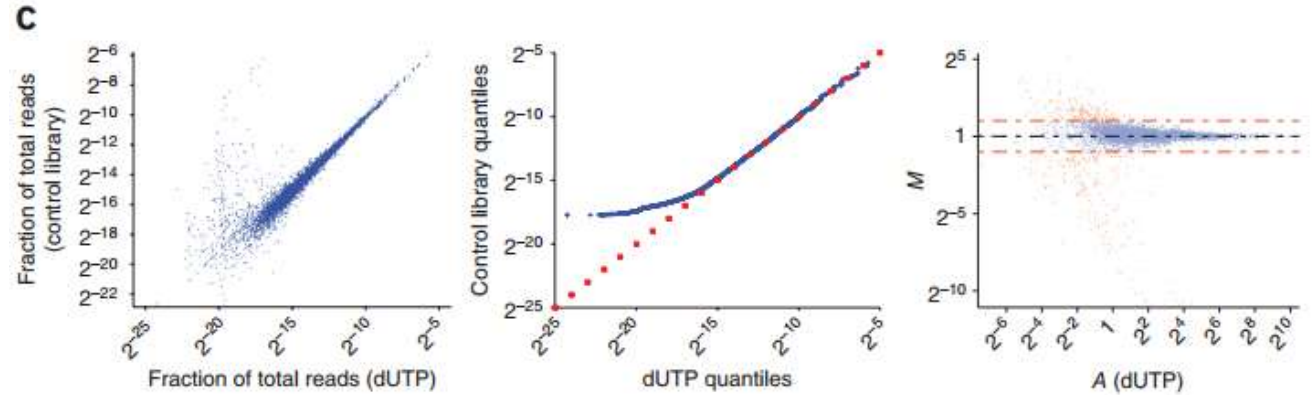
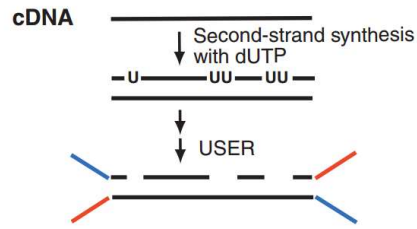
## Wady:

- Stosunkowo skomplikowana procedura

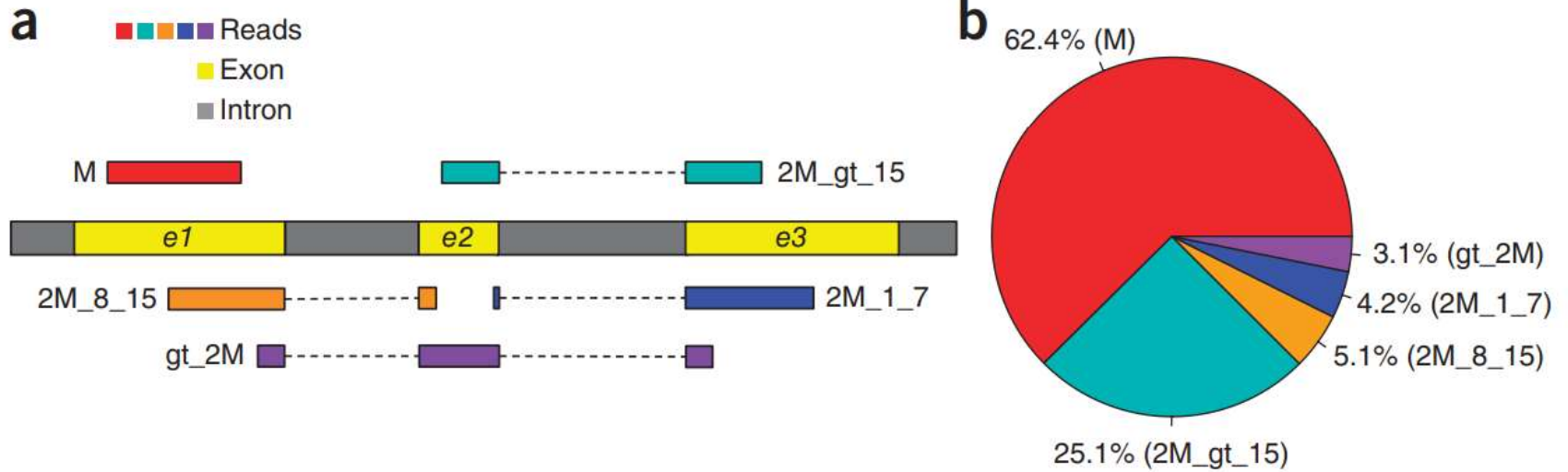
## Zastosowanie:

- Badanie poziomu ekspresji GENÓW i TRANSKRYPTÓW
- Badanie aktywności transkrypcyjnej genomu
- Badanie struktury genów
- Identyfikacja izoform splicingowych
- Badanie polimorfizmu nukleotydowego na poziomie transkryptomu

# profilowanie ekspresji genów - RNA-seq dUTP



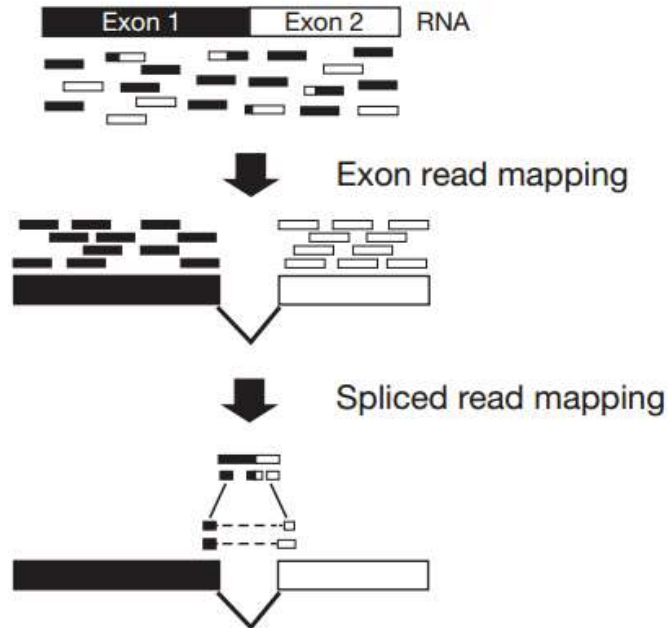
# dopasowanie z przerwami?



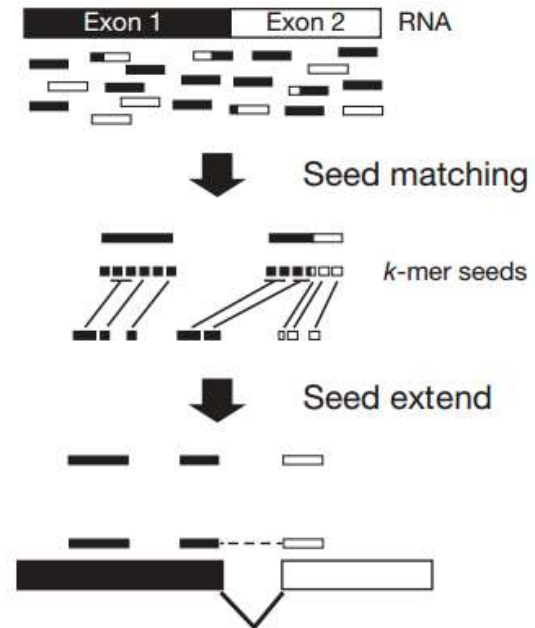
Tak! Dla RNA-Seq!

# dopasowanie z przerwami

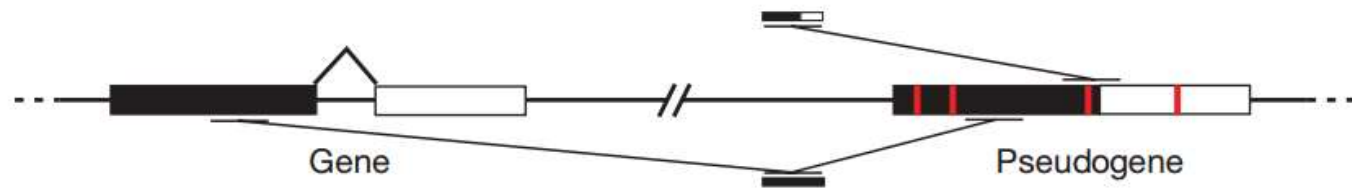
## a Exon-first approach



## b Seed-extend approach

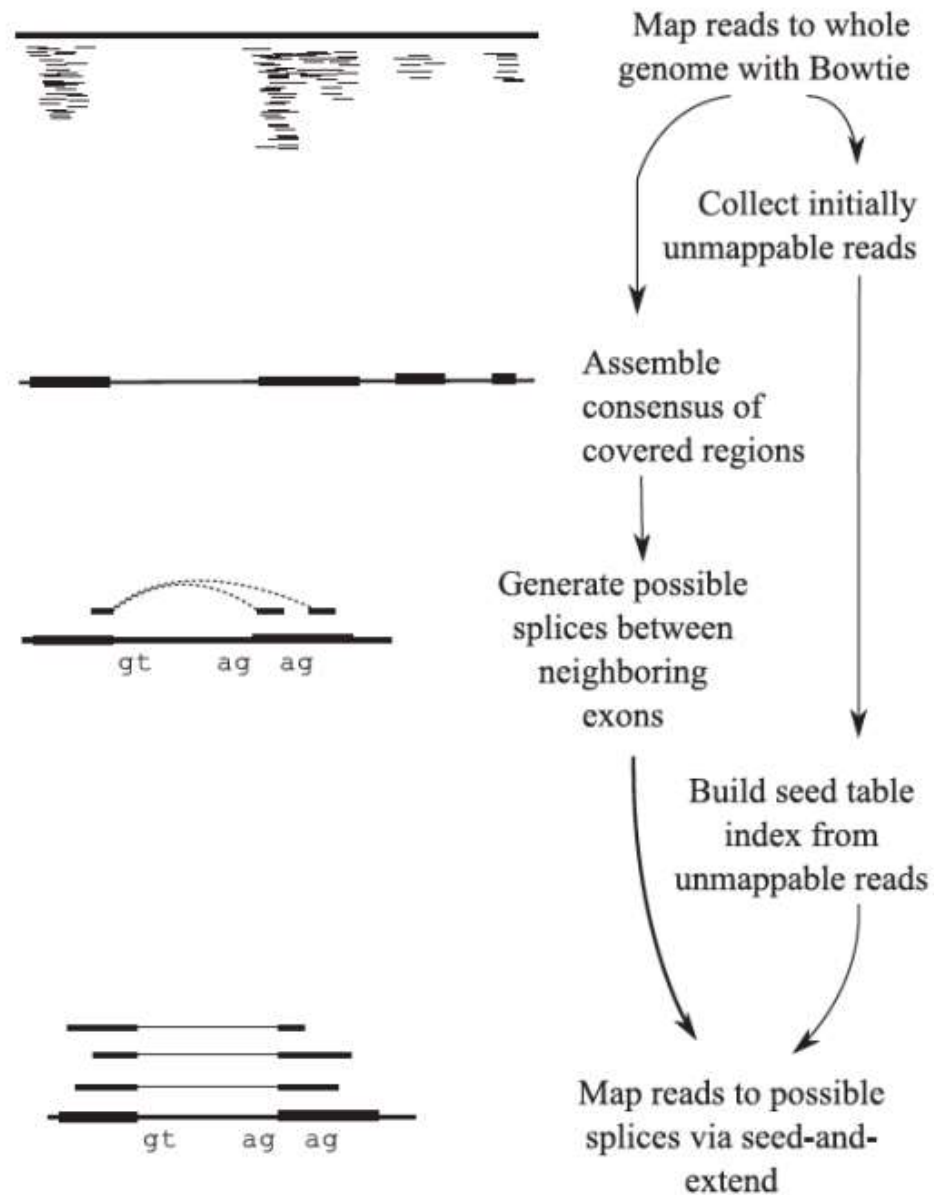


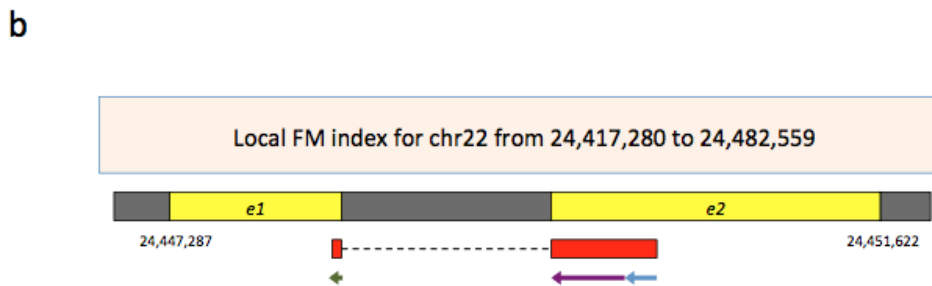
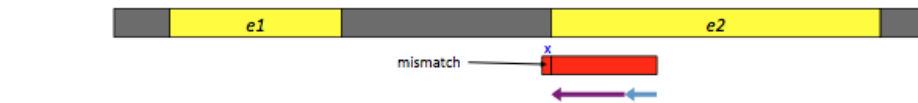
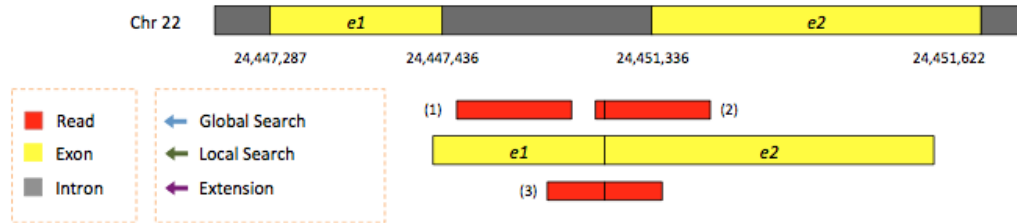
## c Potential limitations of exon-first approaches

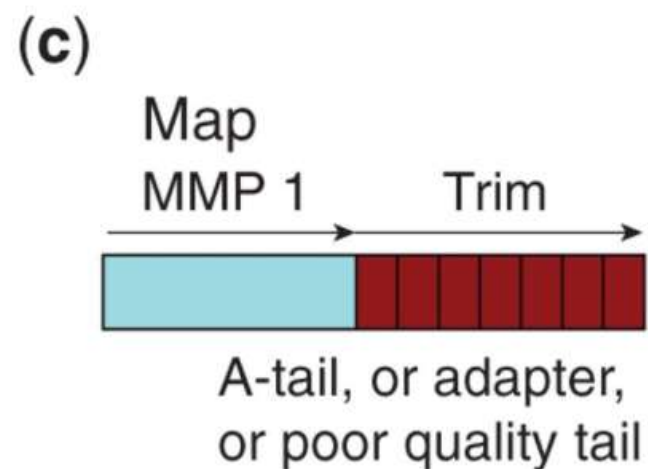
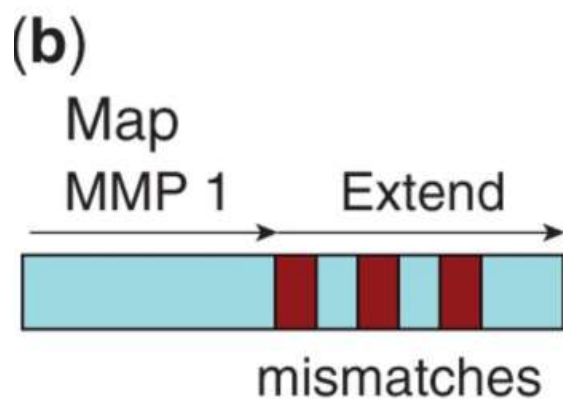
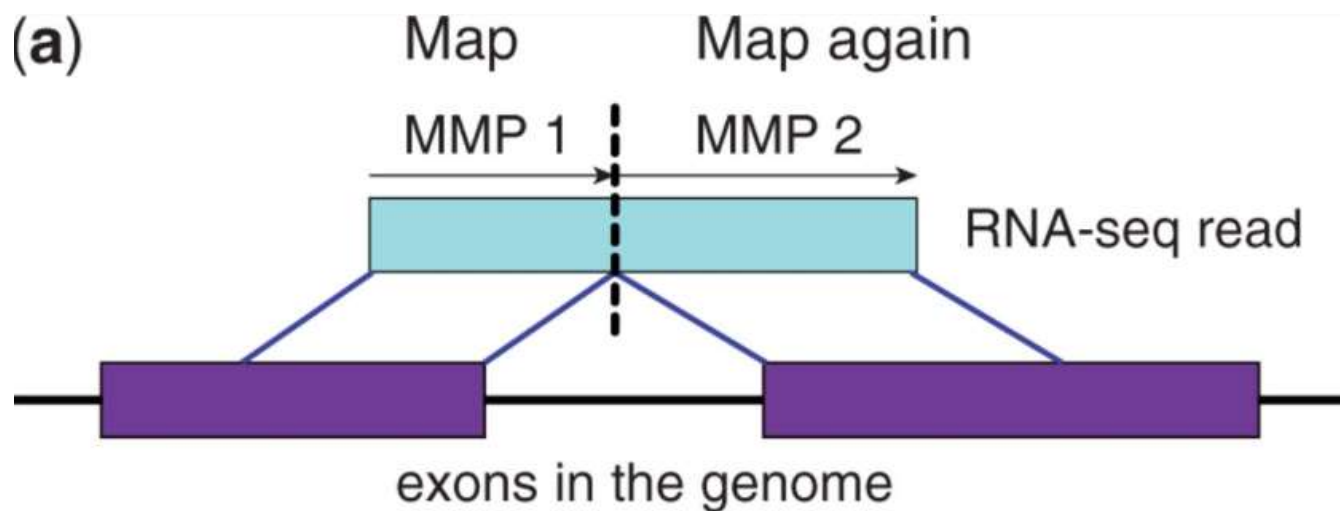




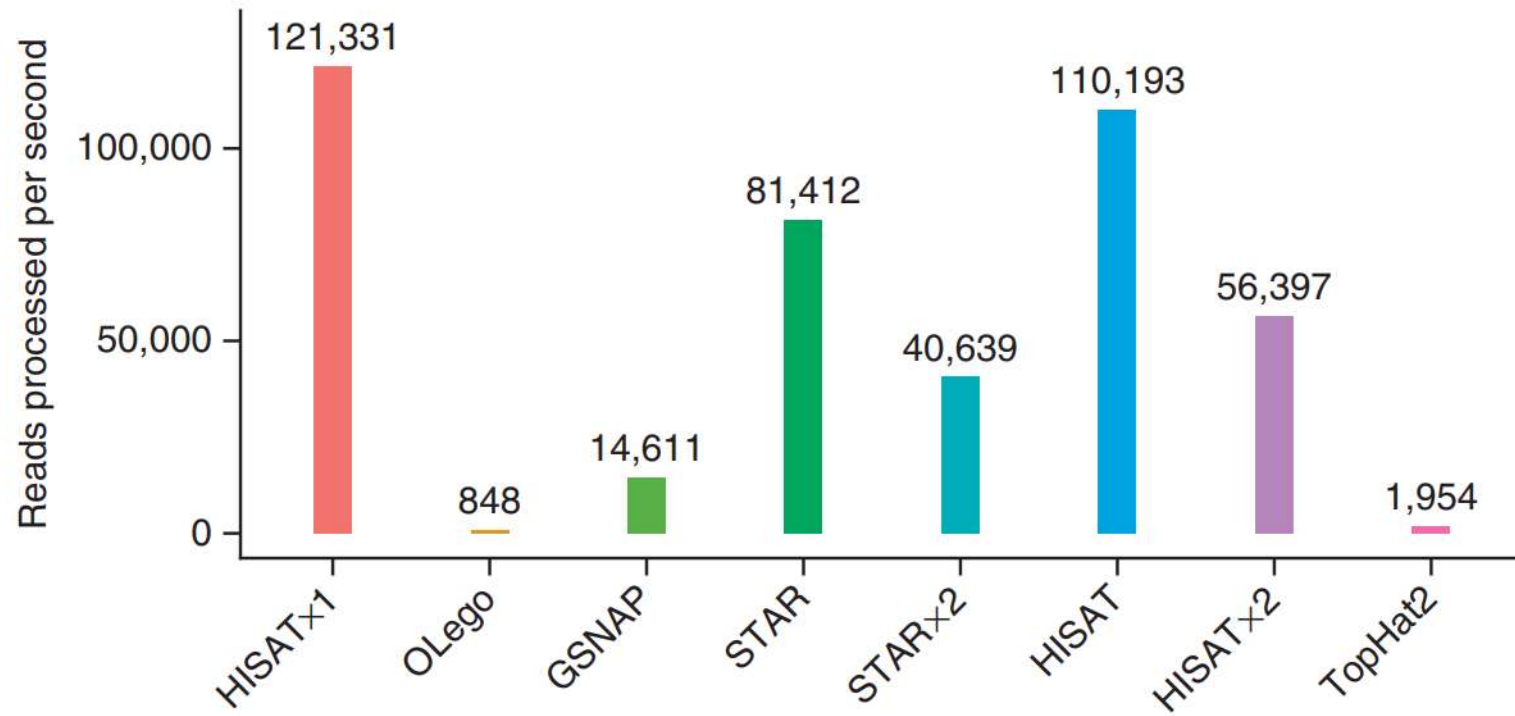
# TopHat



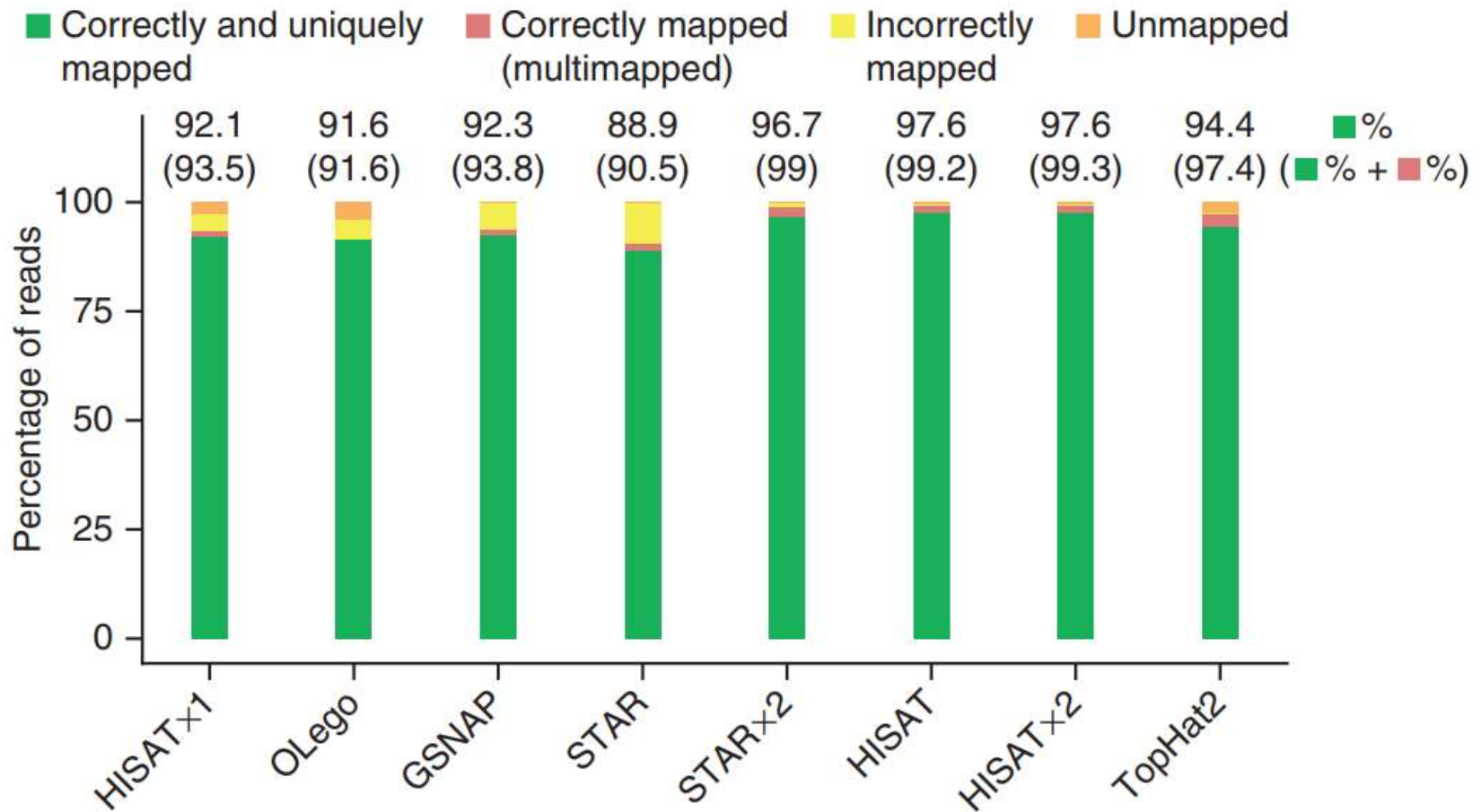




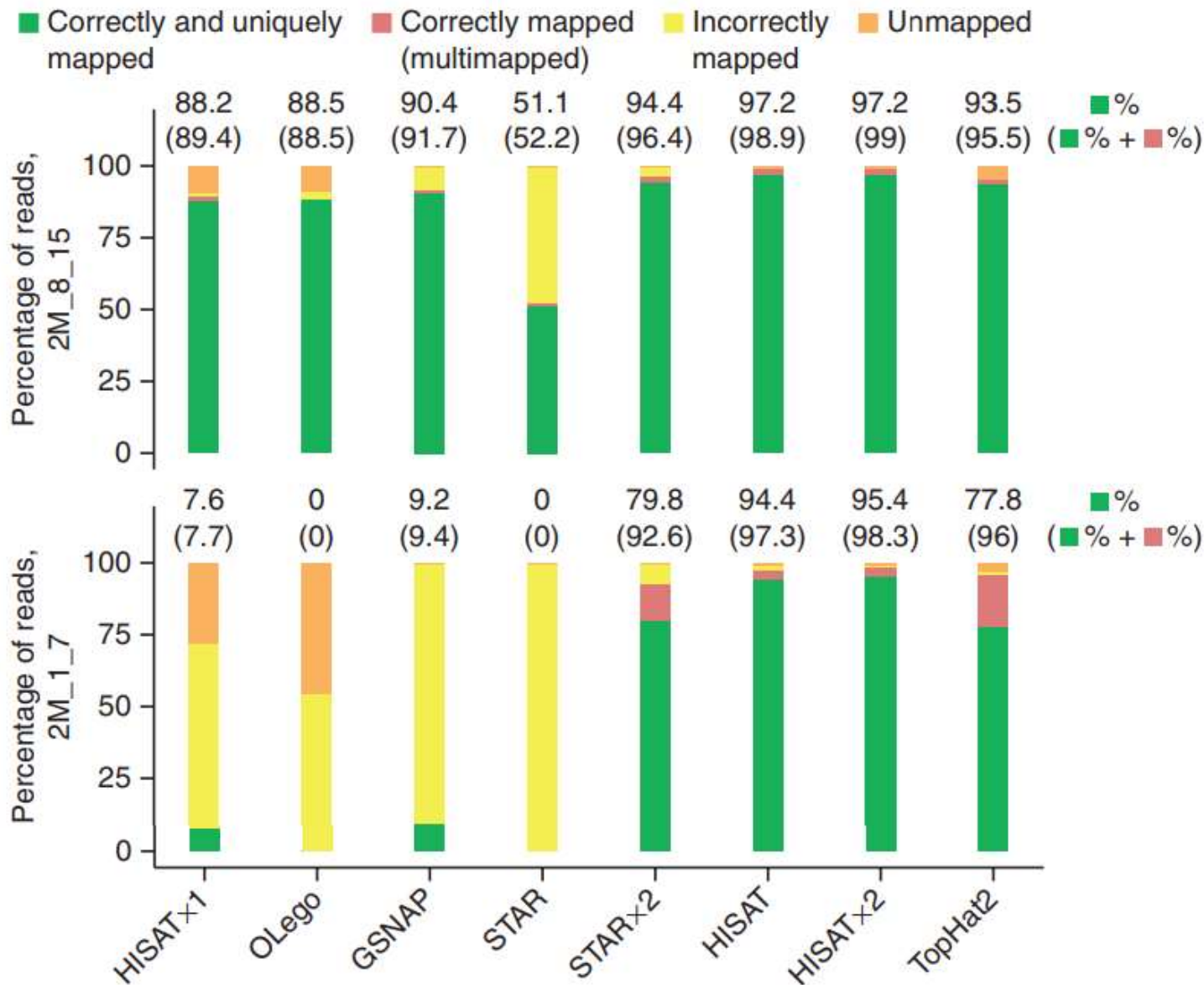
# czas działania



# dokładność

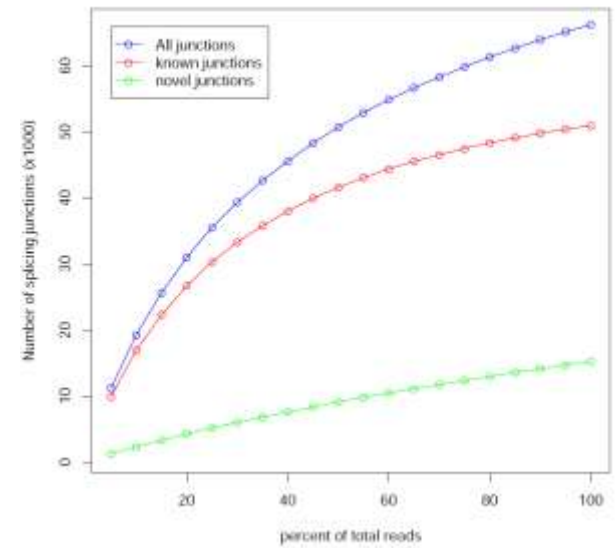
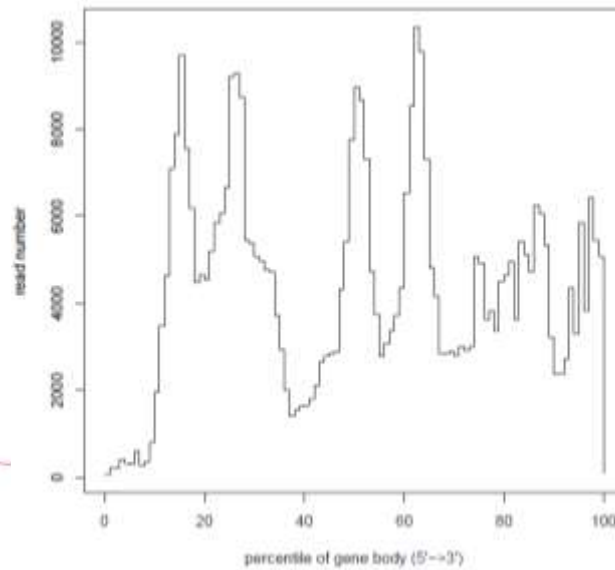
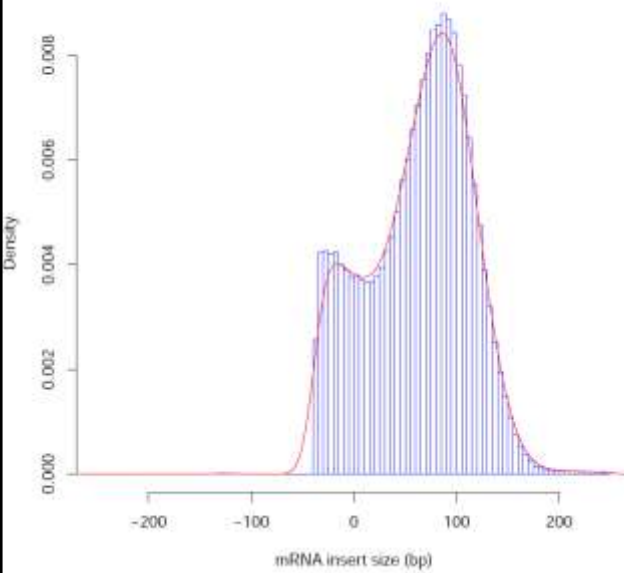


# dokładność



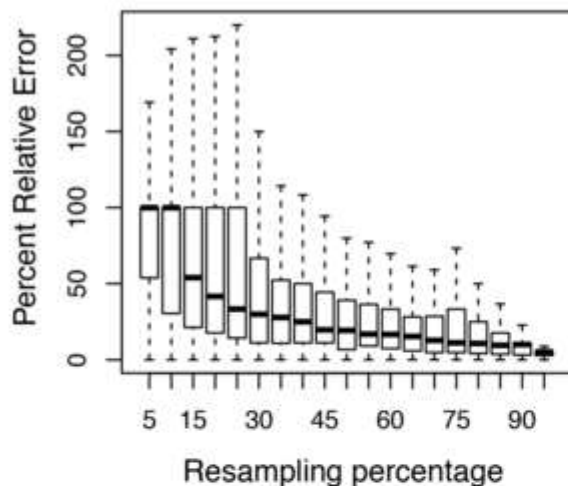
# weryfikacja wyników mapowania: RseQC

Mean=61.7058439571647; SD=50.3152604799405

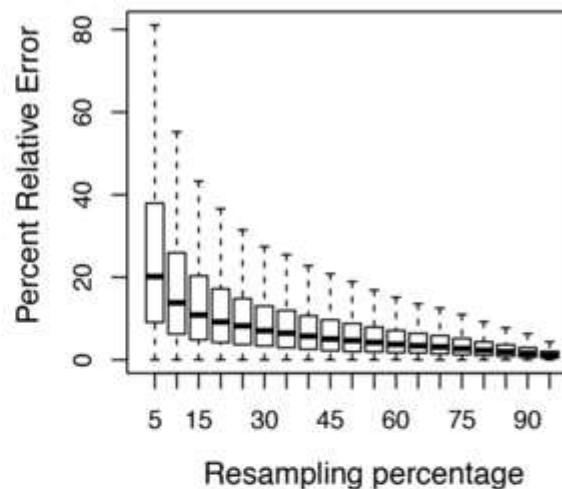


# weryfikacja wyników mapowania: RseQC

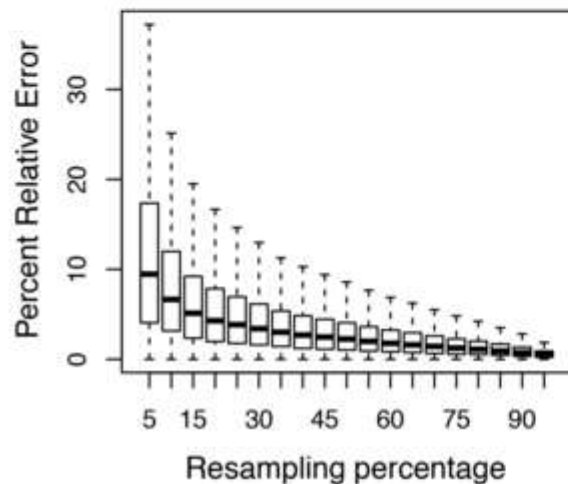
Q1



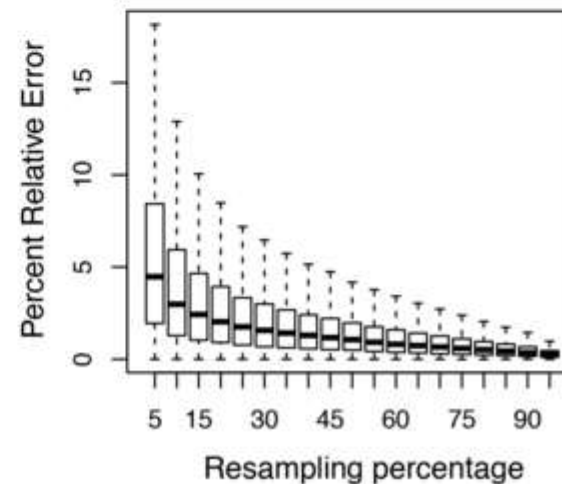
Q2



Q3



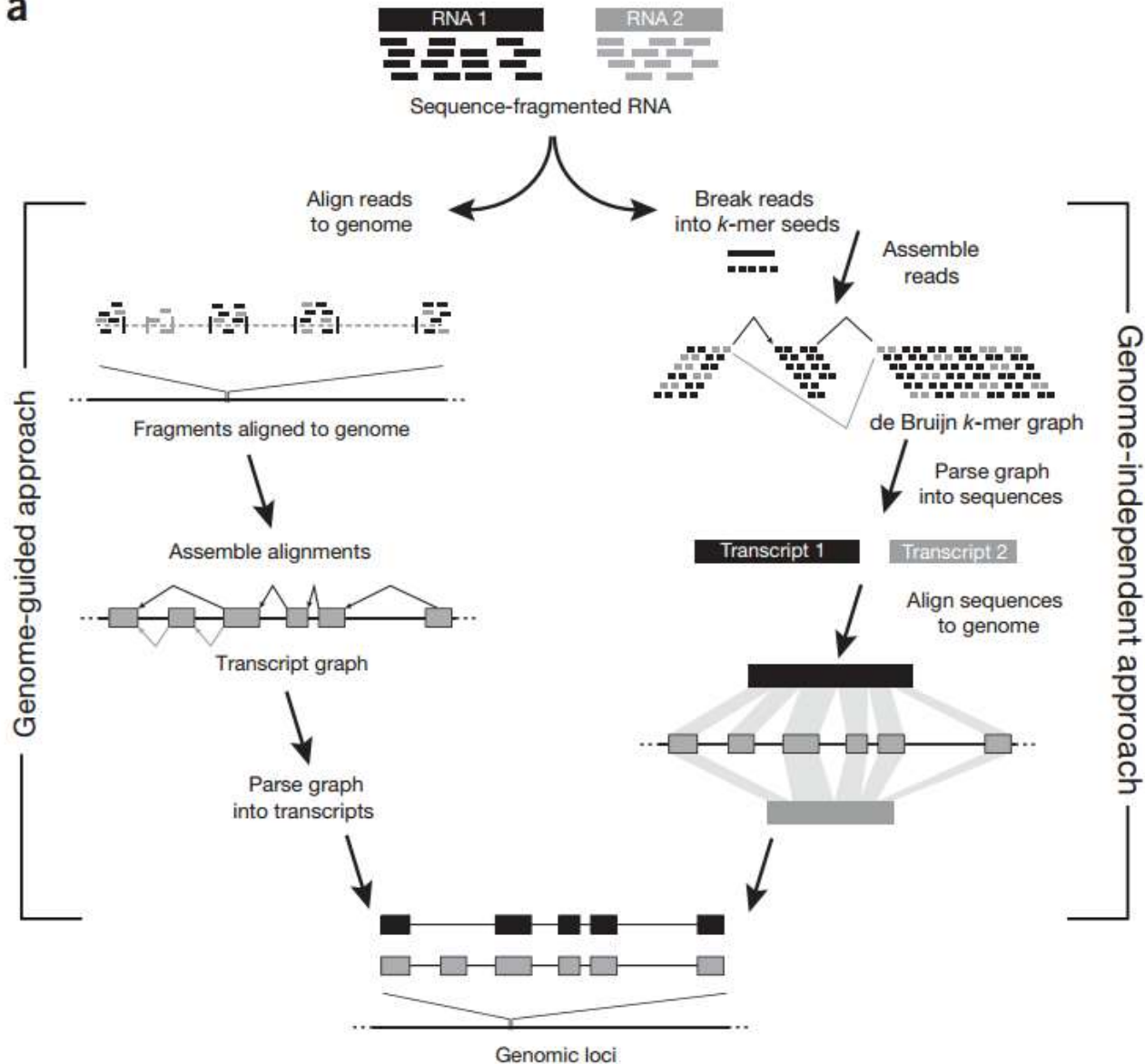
Q4





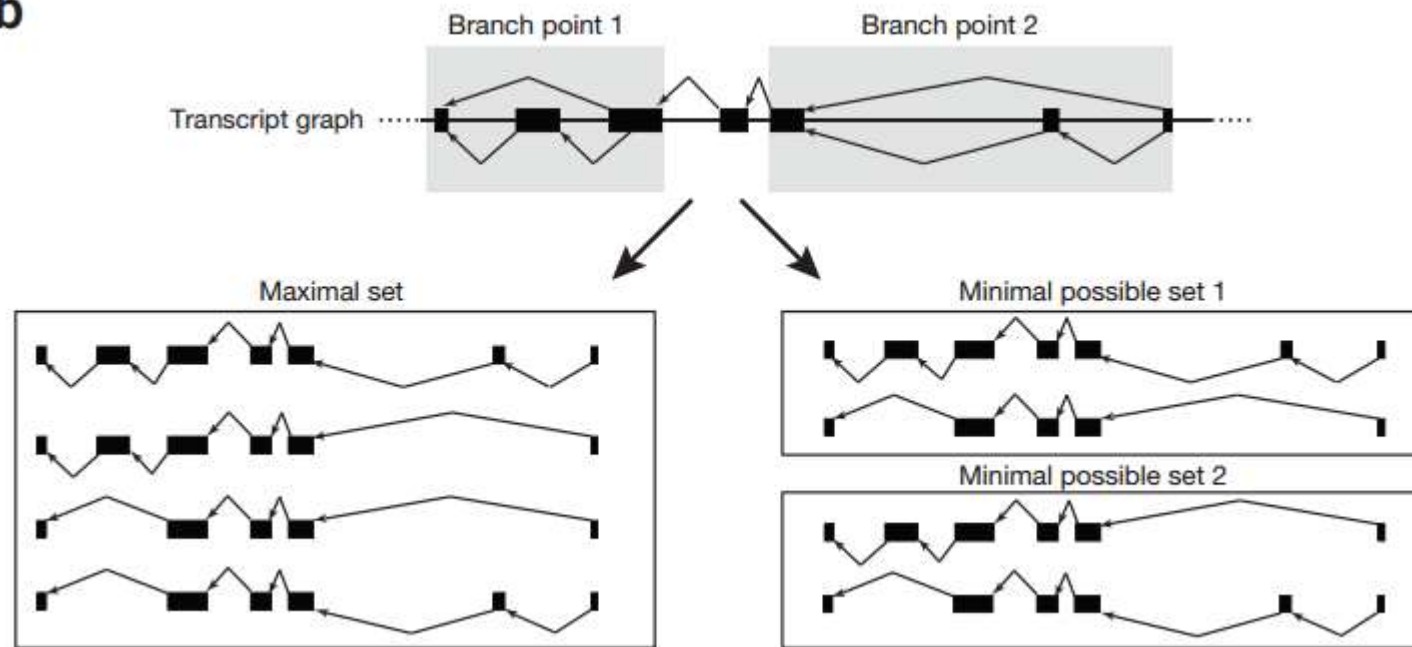
# składanie transkryptów

a



# składanie transkryptów

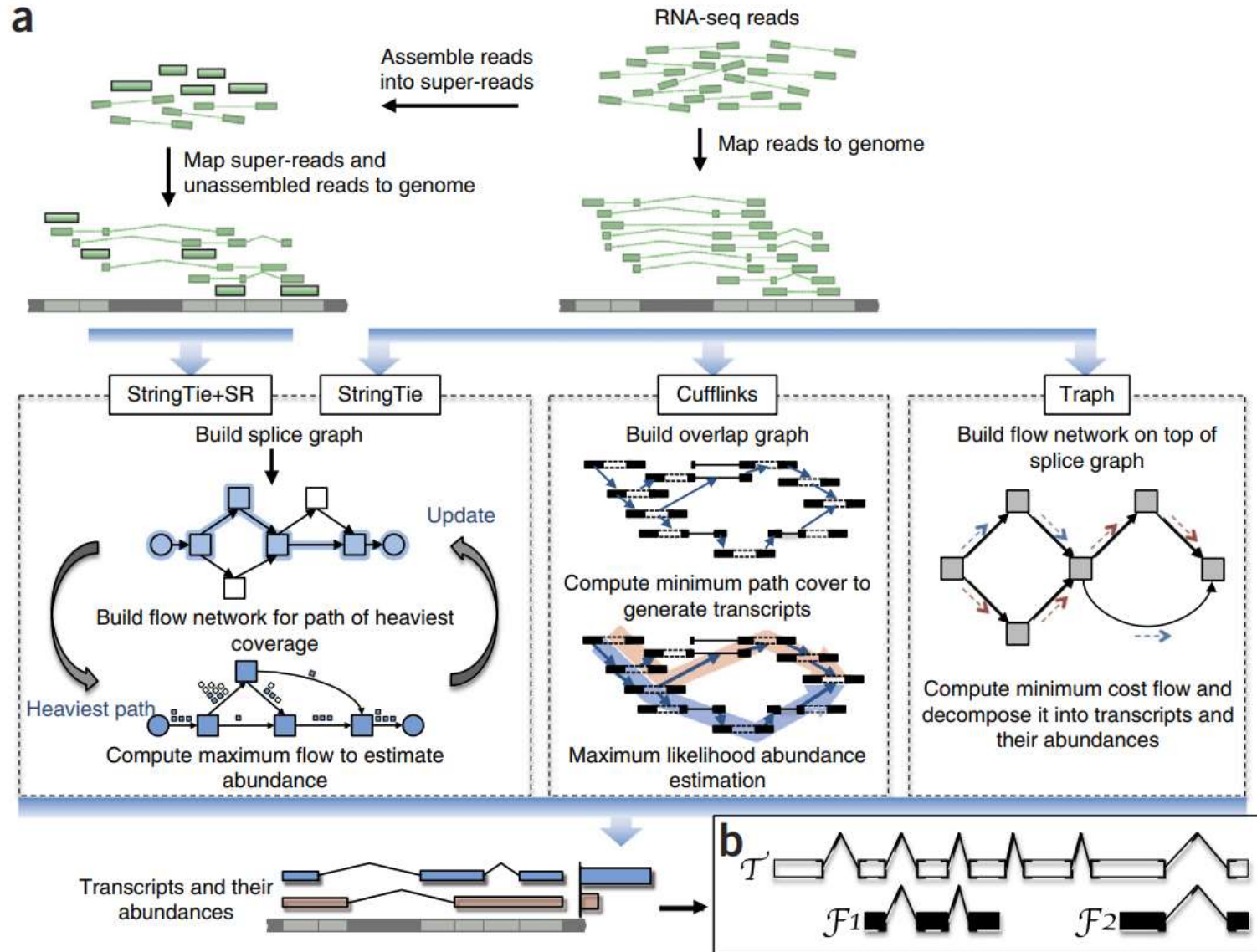
**b**



sensitivity  
(Scripture)

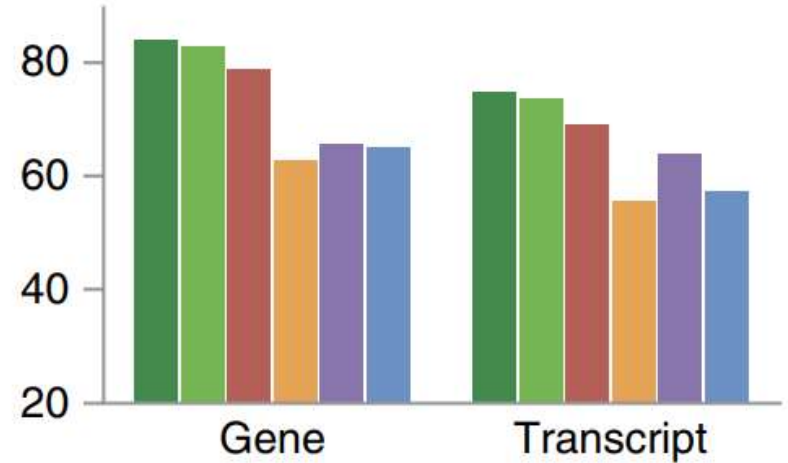
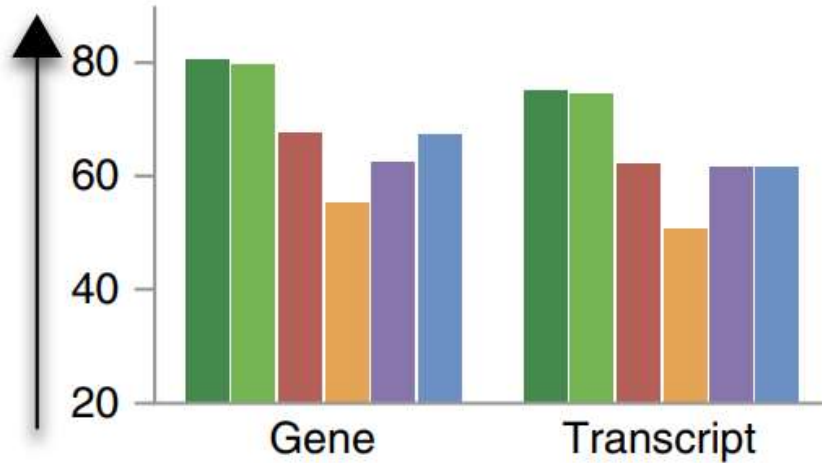
precision  
(Cufflinks)

# składanie transkryptów

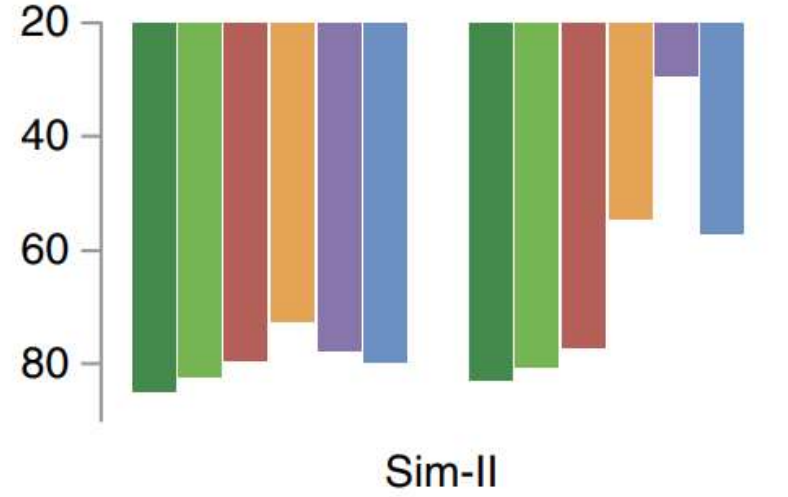
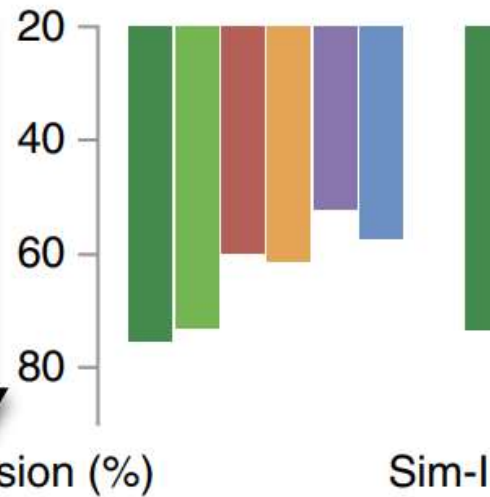


# składanie transkryptów

Sensitivity (%)



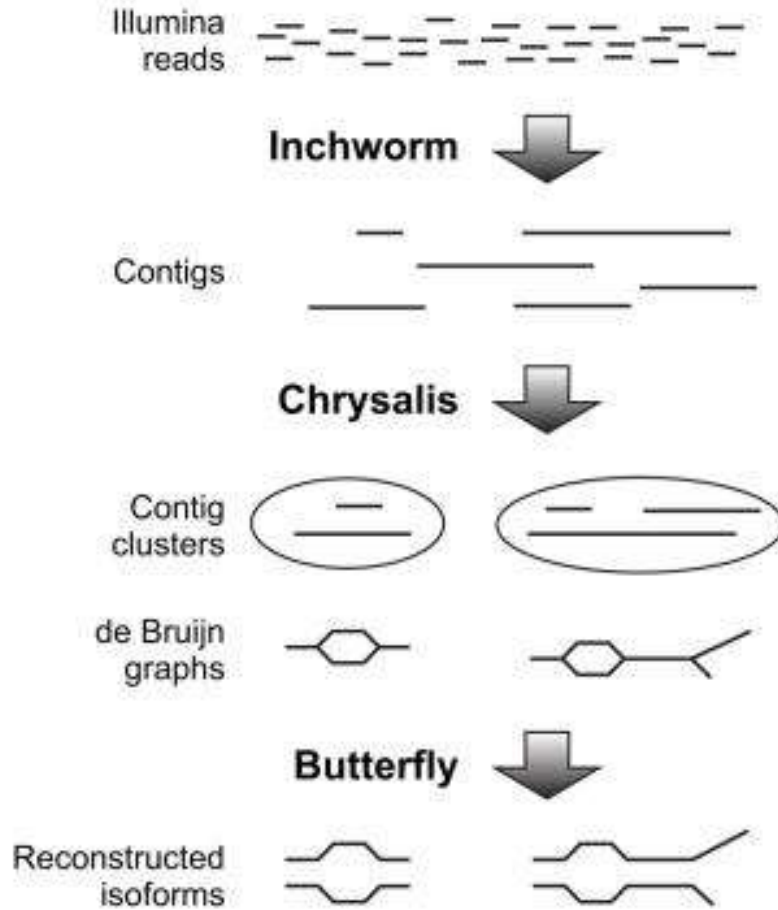
Precision (%)



StringTie+SR StringTie Cufflinks Traph Scripture IsoLasso

# składanie transkryptomu *de novo*

## Trinity



Indication of compute resources

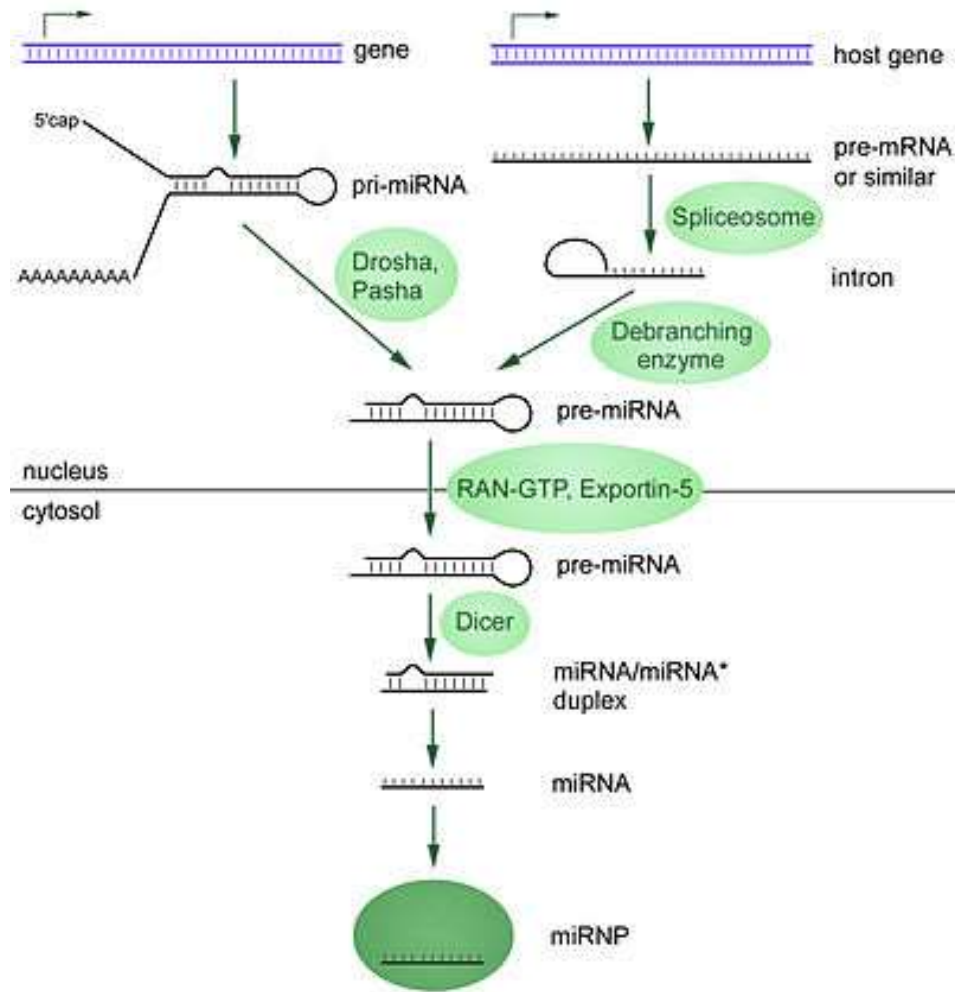


Single high-memory, multi-core server

Massively parallel on computing grid



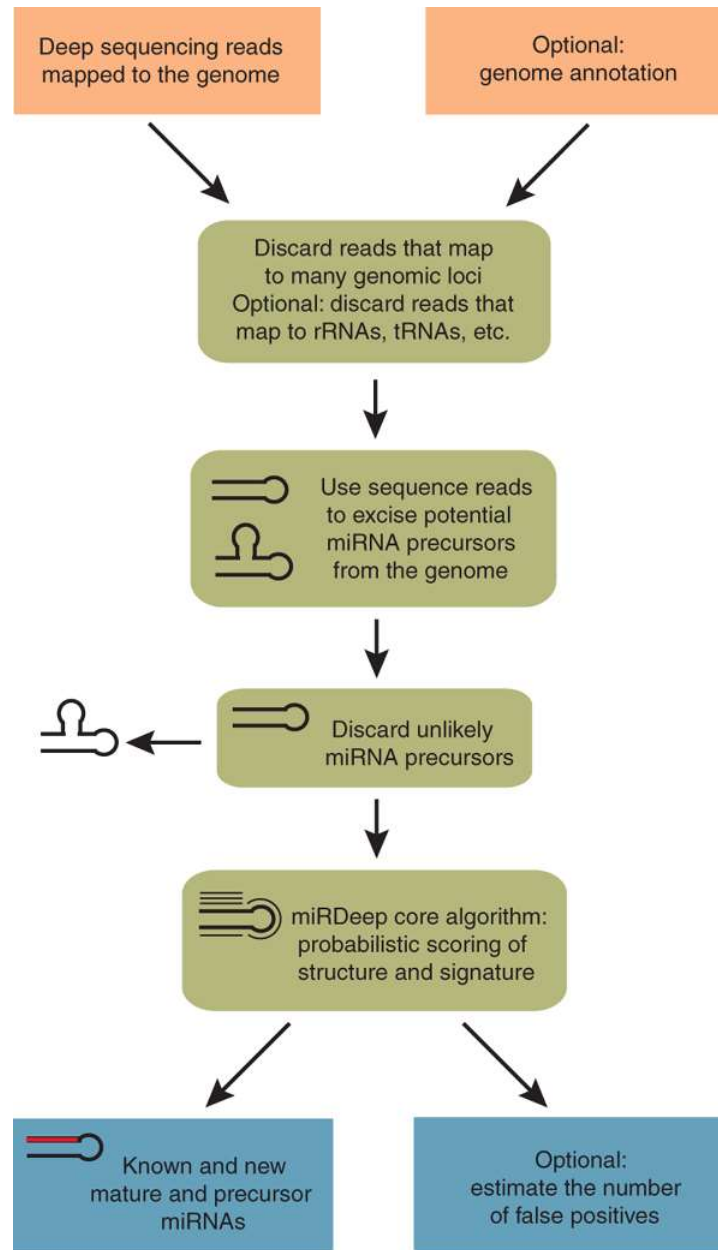
# small RNA-seq



# small RNA-seq

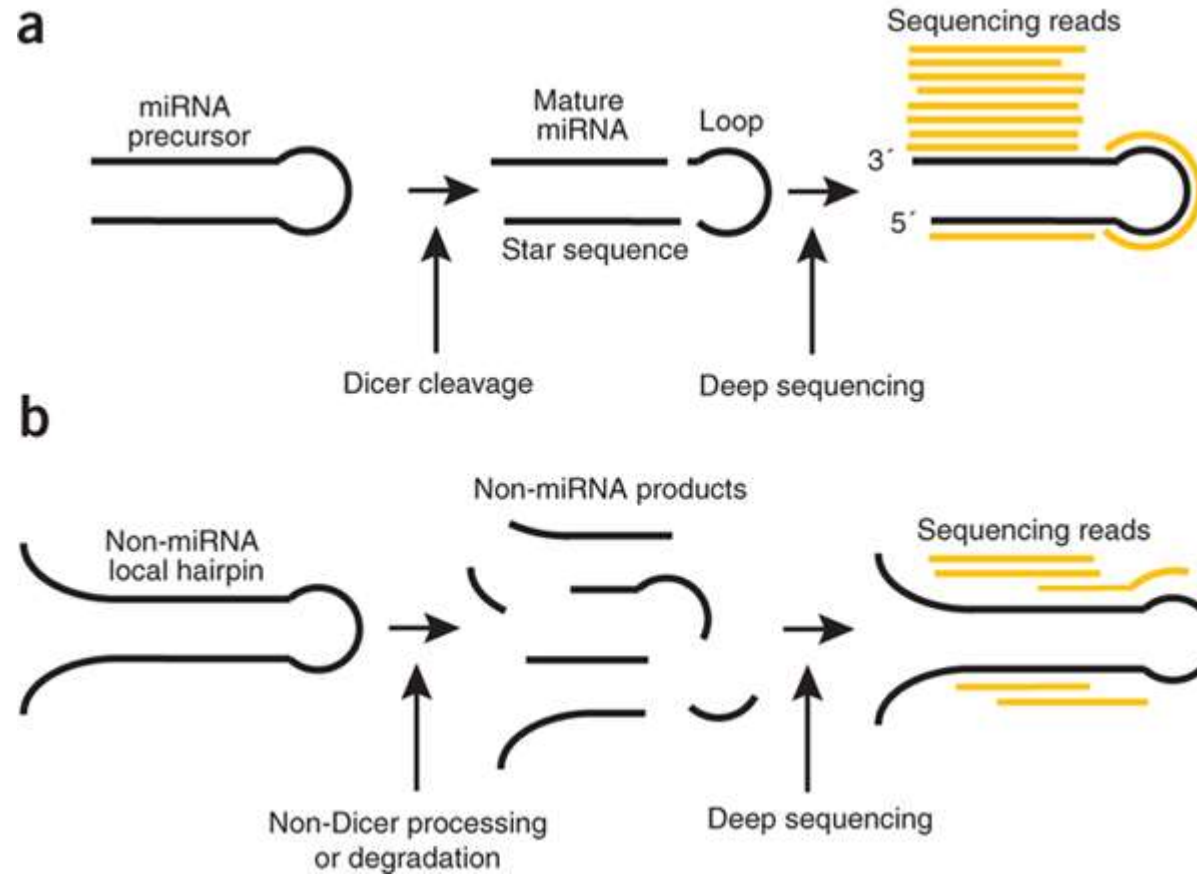
1. Izolacja małych RNA (18-30nt)
2. Brak fragmentacji
3. Sekwencjonowanie pojedynczego końca
4. Mapowanie bez uwzględnienia splicingu

# small RNA-seq mirDeep

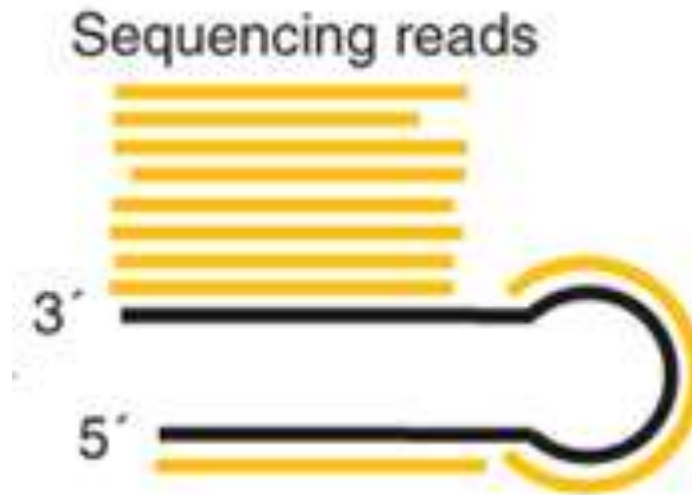




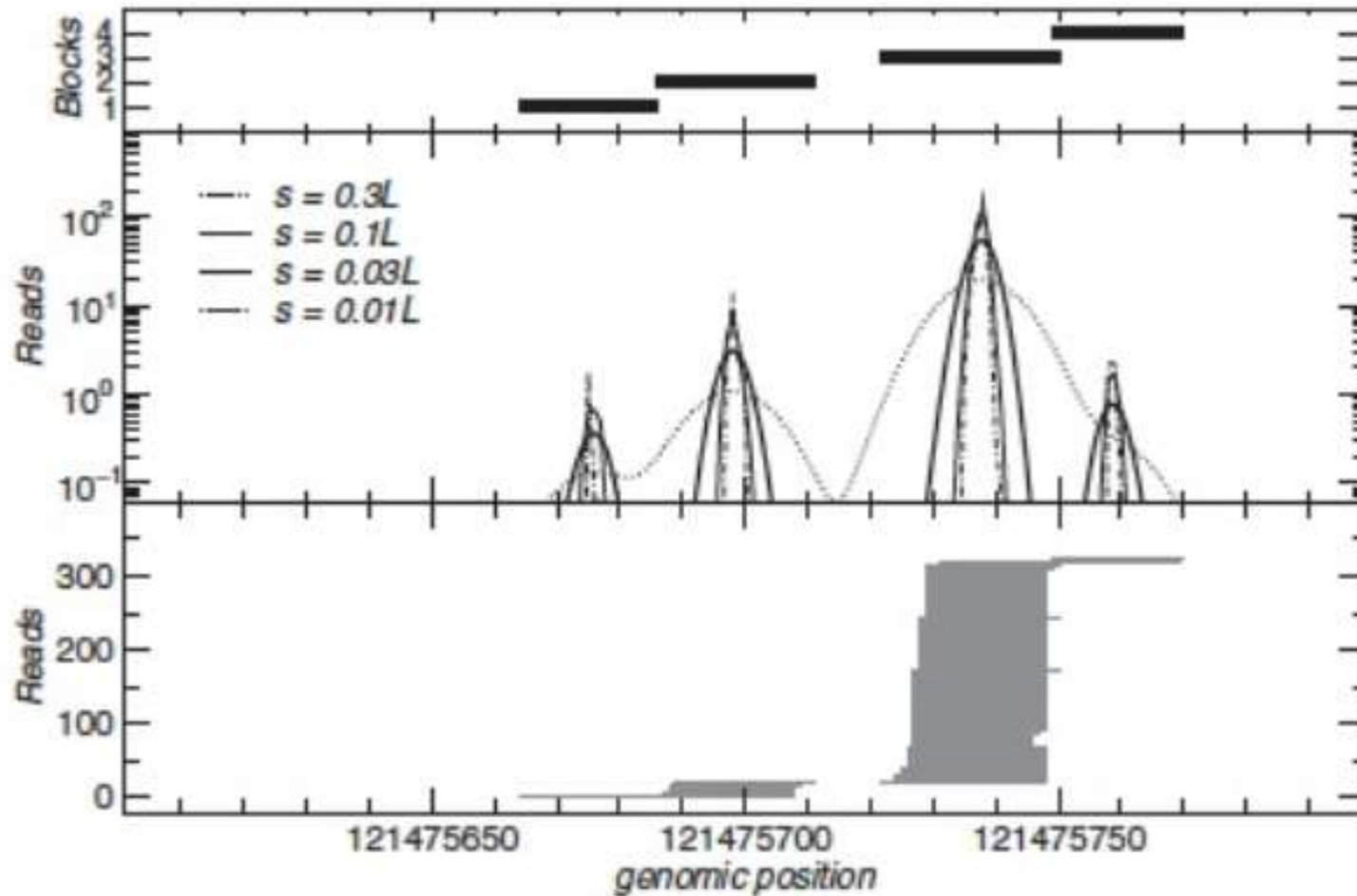
# small RNA-seq mirDeep



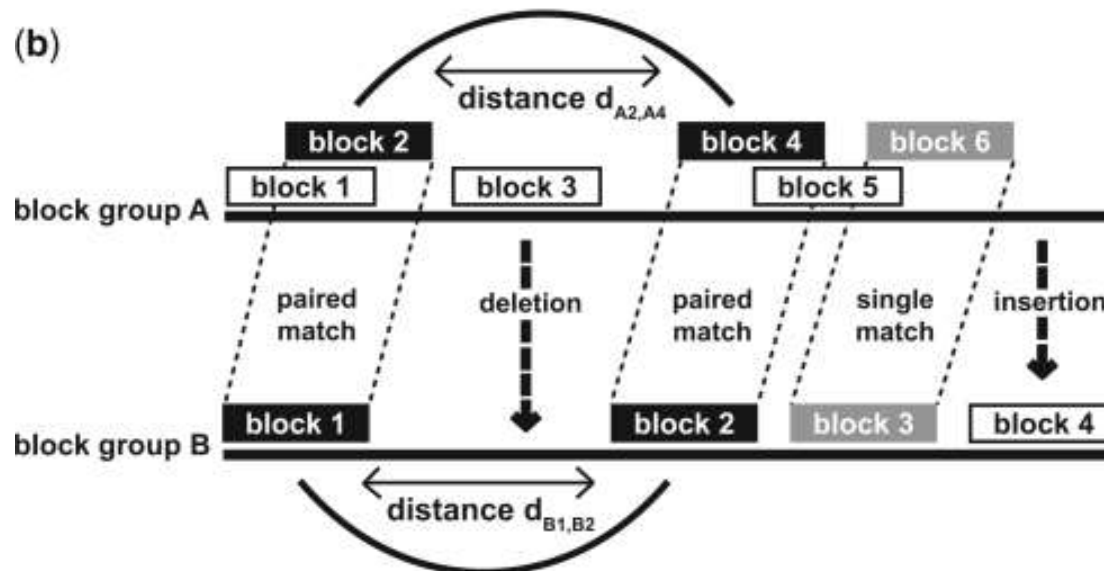
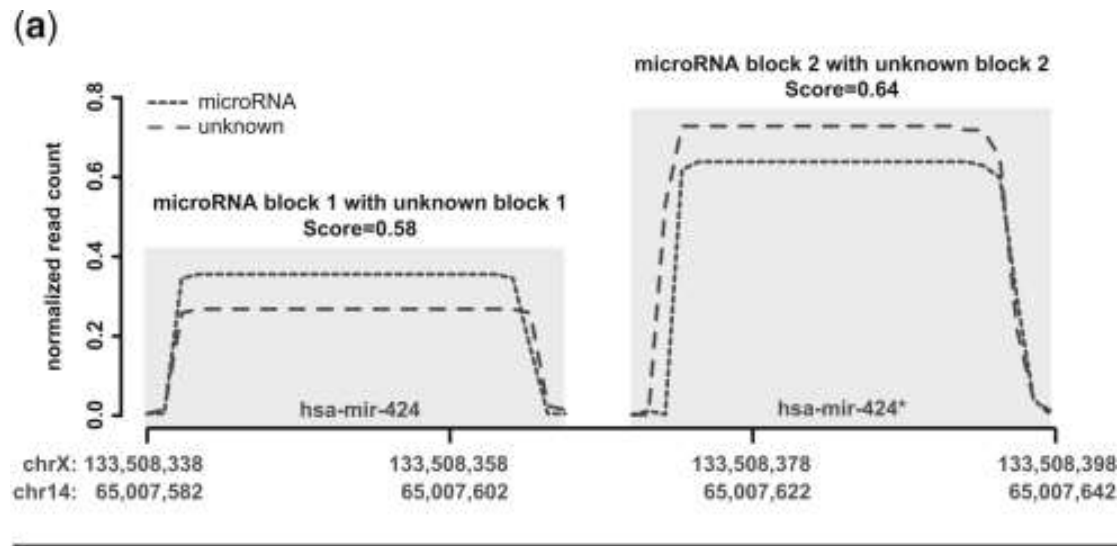
# small RNA-seq poziom ekspresji a duplikaty



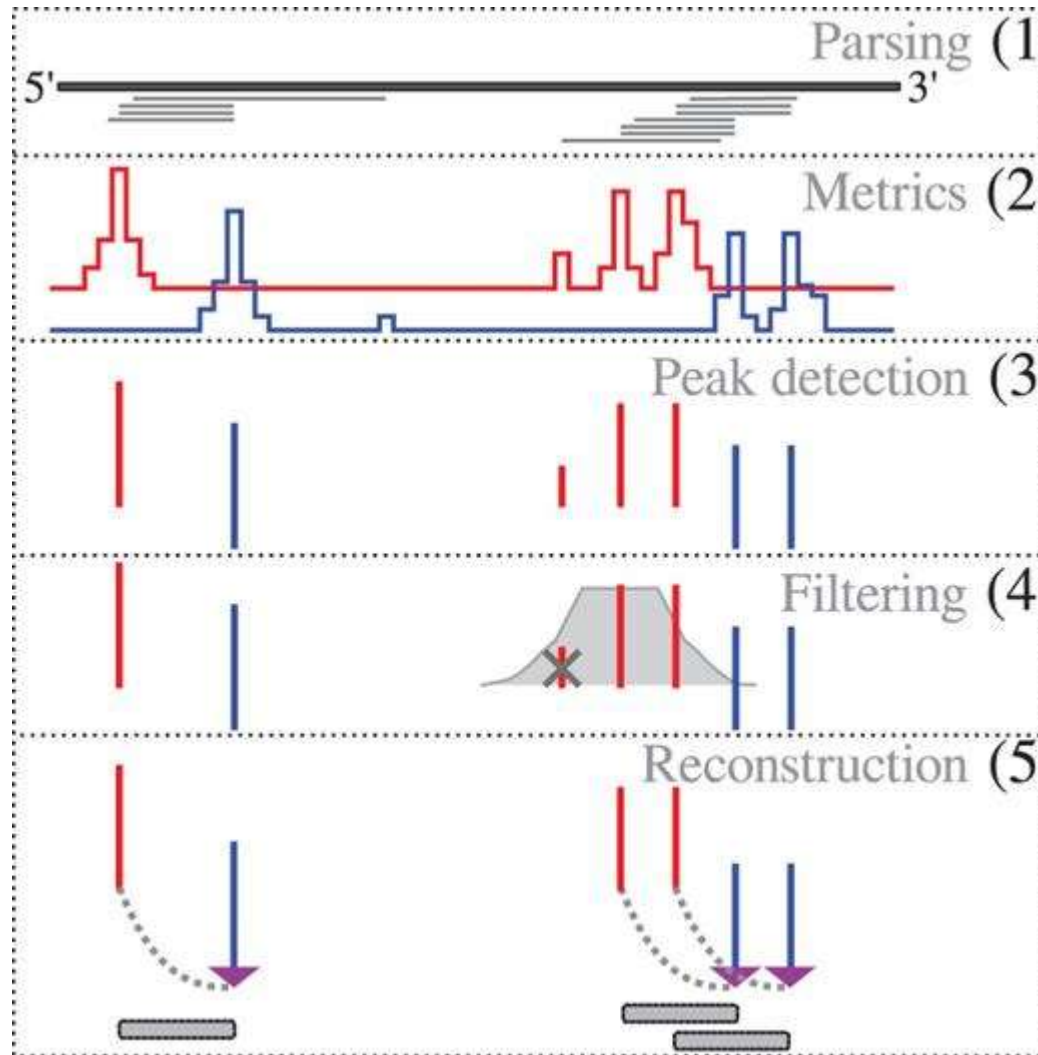
# small RNA-seq blockbuster



# small RNA-seq deepBlockAlign



# small RNA-seq FlaiMapper



# zliczanie poziomu ekspresji

Zliczanie na podstawie adnotacji odczytów przypadających na:

- gen
- transkrypt
- egzon
- ...

Programy:

HTSeq-count

- popularny
- bardzo wolny
- zwraca również geny nie ulegające ekspresji

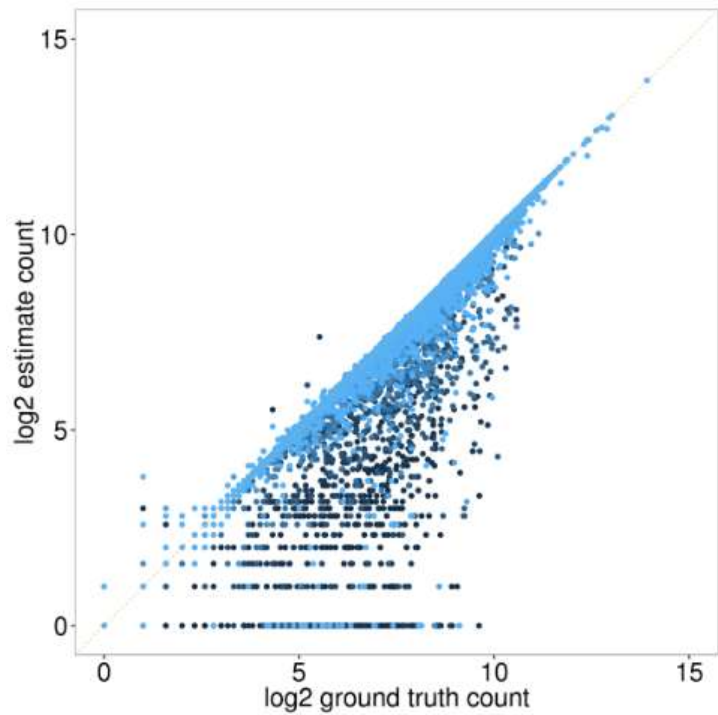
FeatureCounts:

- szybki
- analiza wielu próbek jednocześnie
- duża kontrola nad sposobem zliczania

# metody alignment-free

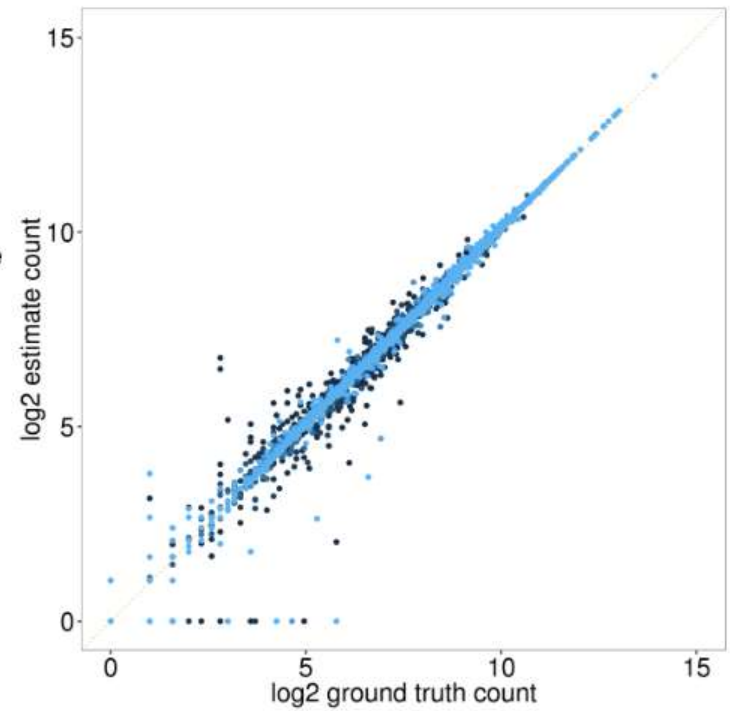
- Brak konieczności mapowania odczytów
- Dopasowanie odczytów na podstawie dystansu (np. częstotliwości k-merów)
- Dokładność niezależna od „mapowalności” odczytów
- Bardzo szybkie
- Wyniki problematyczne w interpretacji podczas testowania różnicowej ekspresji
  
- Kallisto
- Sailfish
- Salomon

# metody alignment-free



Fraction unique

- 1.00
- 0.75
- 0.50
- 0.25
- 0.00

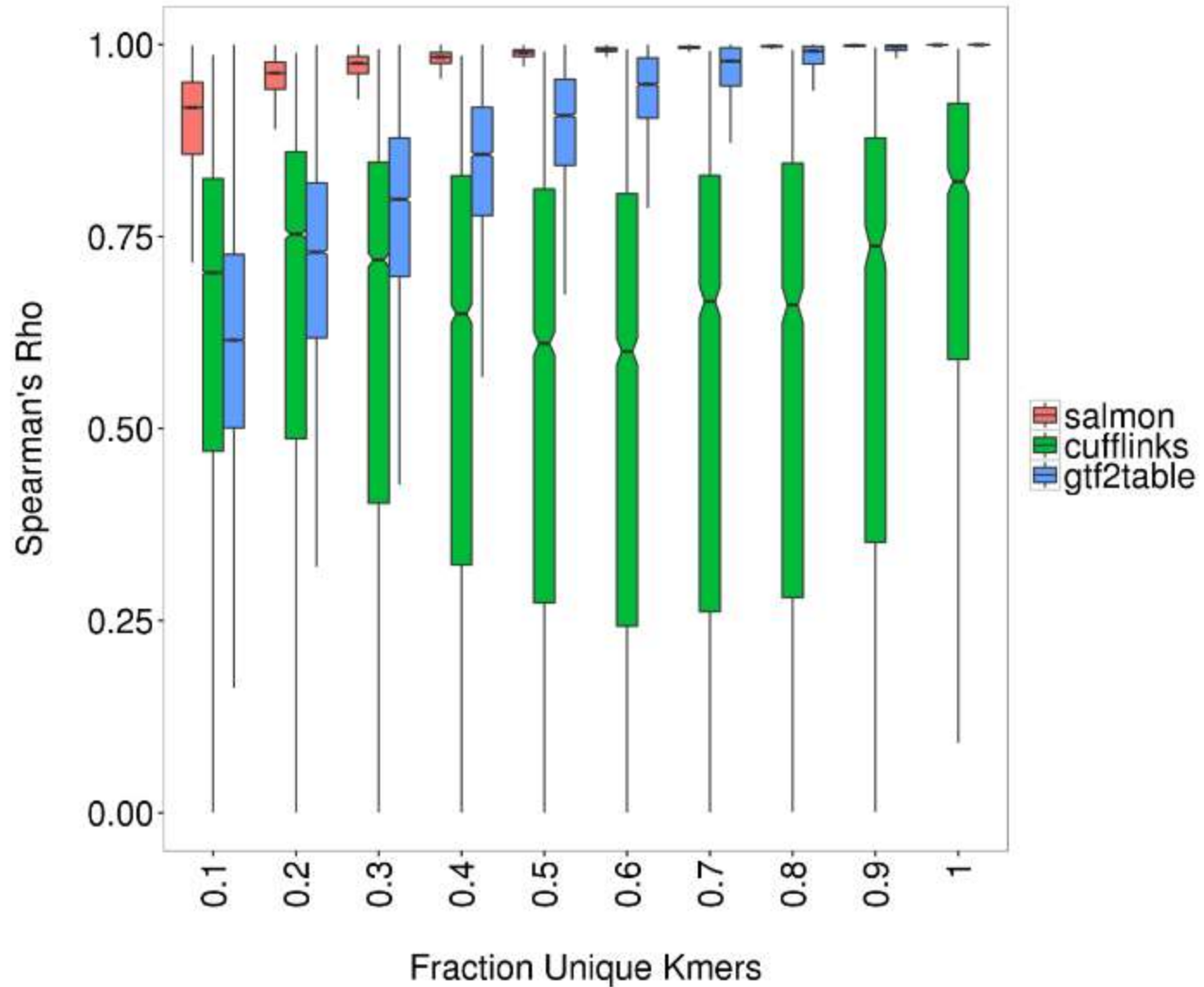


Fraction unique

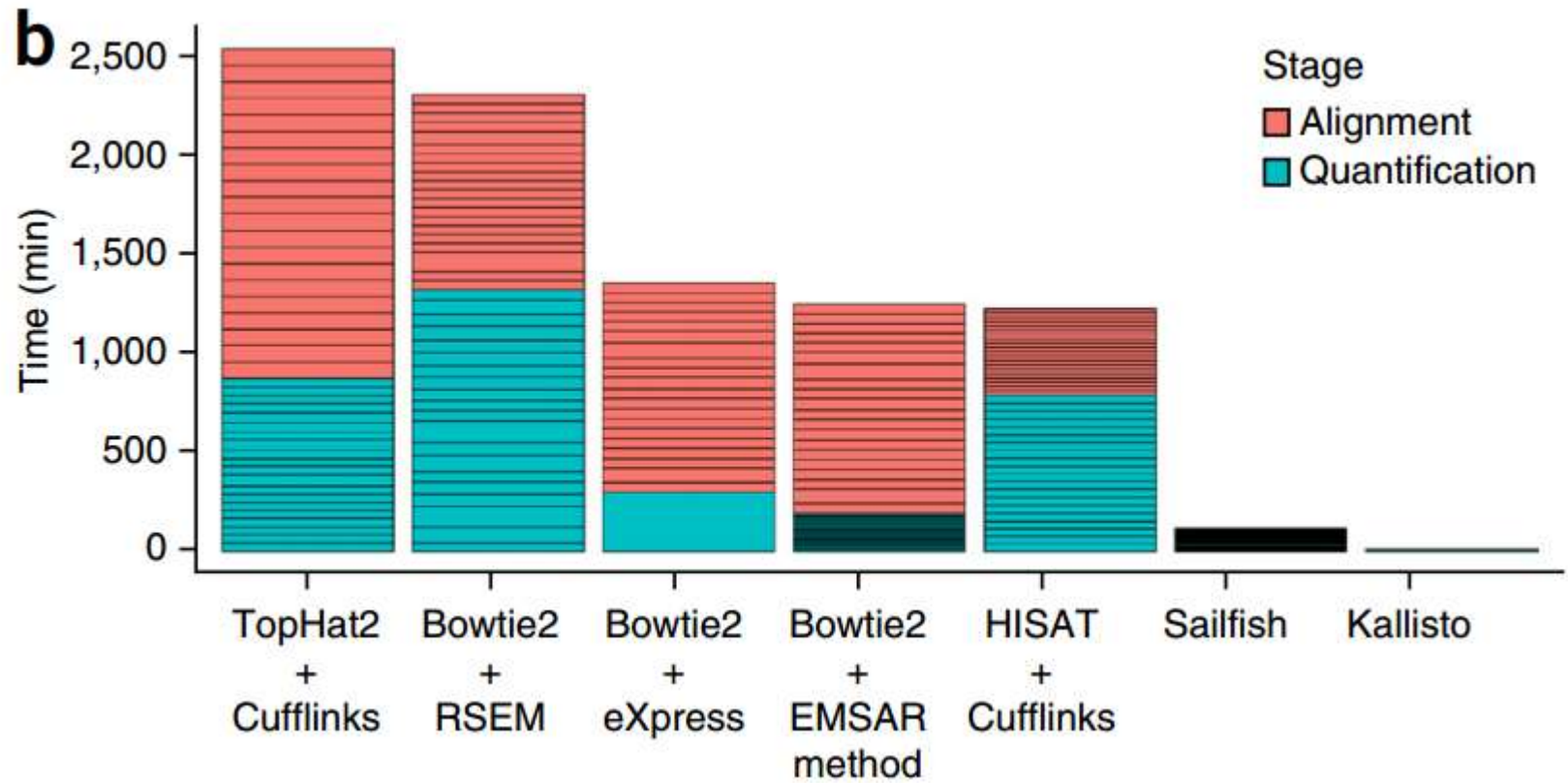
- 1.00
- 0.75
- 0.50
- 0.25
- 0.00



# metody alignment-free



# metody alignment-free



# RNA-seq różnicowa ekspresja

Założenia:

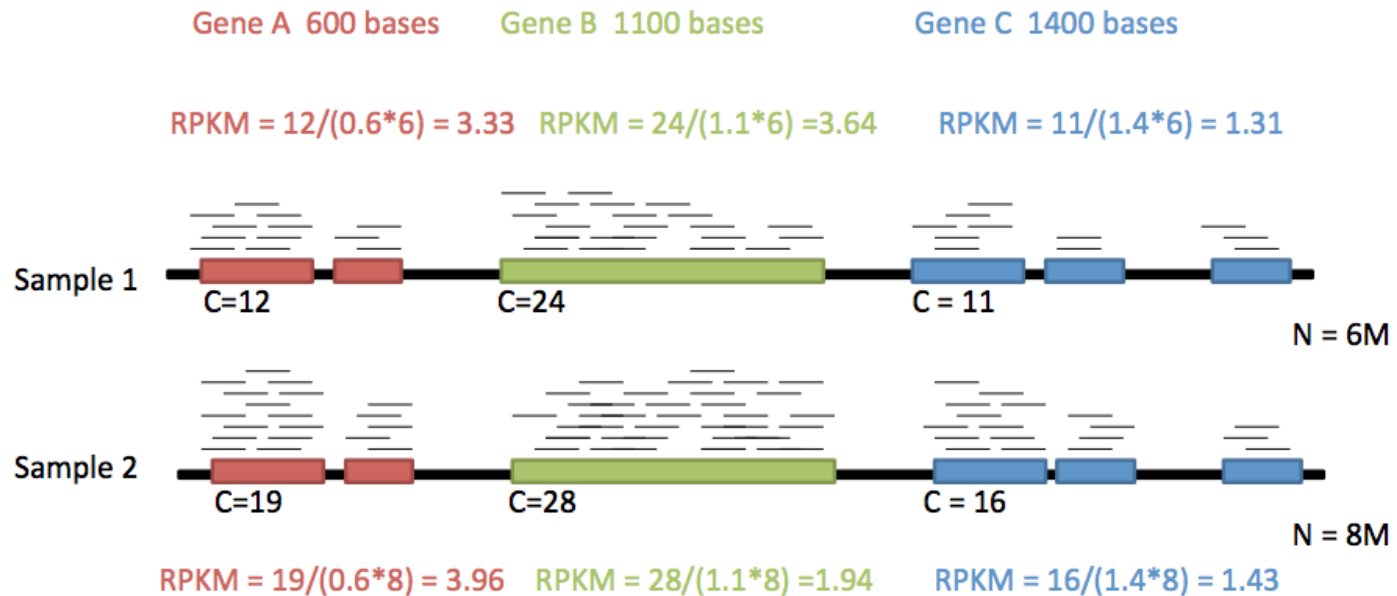
1. Większość genów ulega ekspresji na tym samym poziomie
2. Wariacje w obrębie replikatów technicznych nie występują

Zaburzenia pomiędzy próbkami:

- zróżnicowana ilość odczytów
- zróżnicowana długość genów
- zróżnicowana kompozycja odczytów

# normalizacja poziomów ekspresji

RPKM: Reads Per Kilobase of transcript per Milion of mapped reads



Małe RNA (brak fragmentacji):

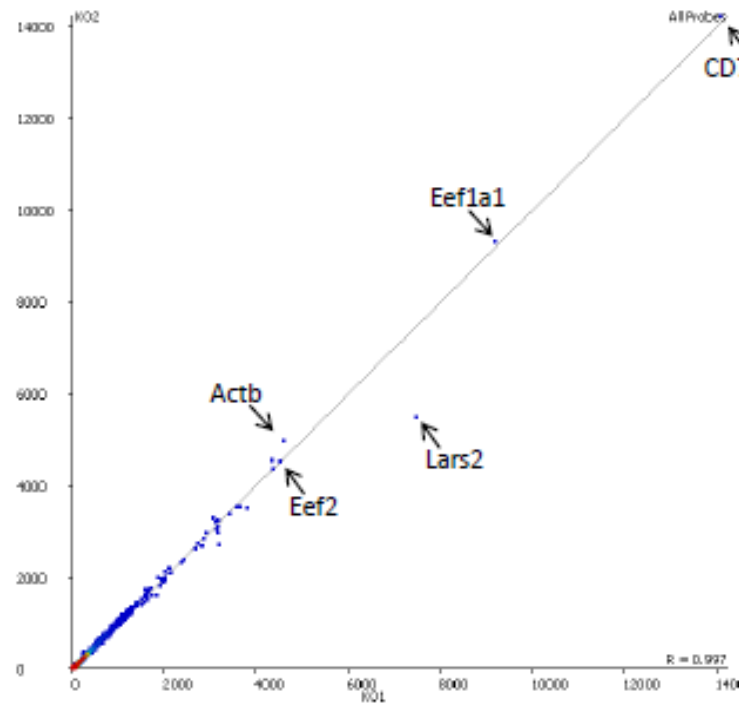
RPM: Reads Per Milion of mapped reads

# wyznaczanie różnic ekspresji

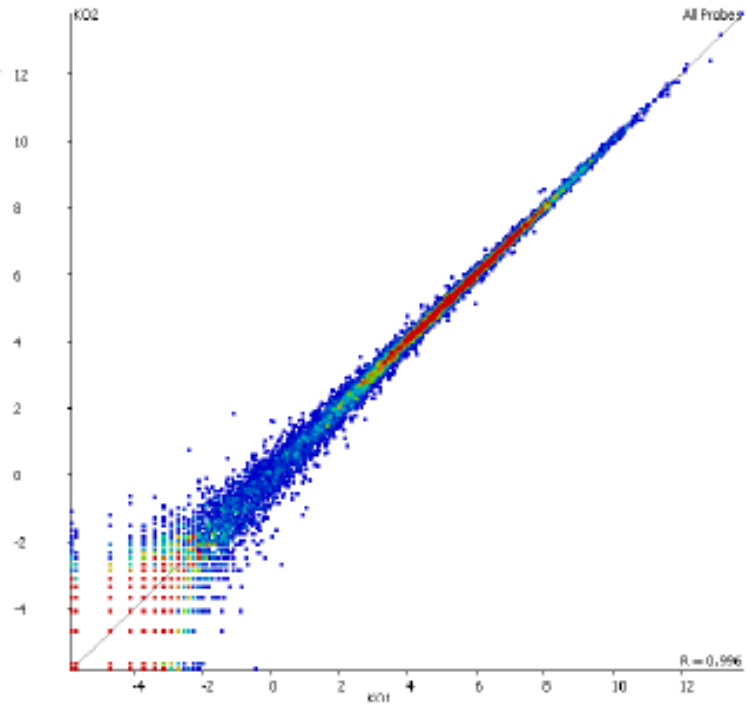
$$\log_2 FC = \log_2(RPKM1/RPKM2)$$

# wyznaczanie różnic ekspresji

Linear



Log2

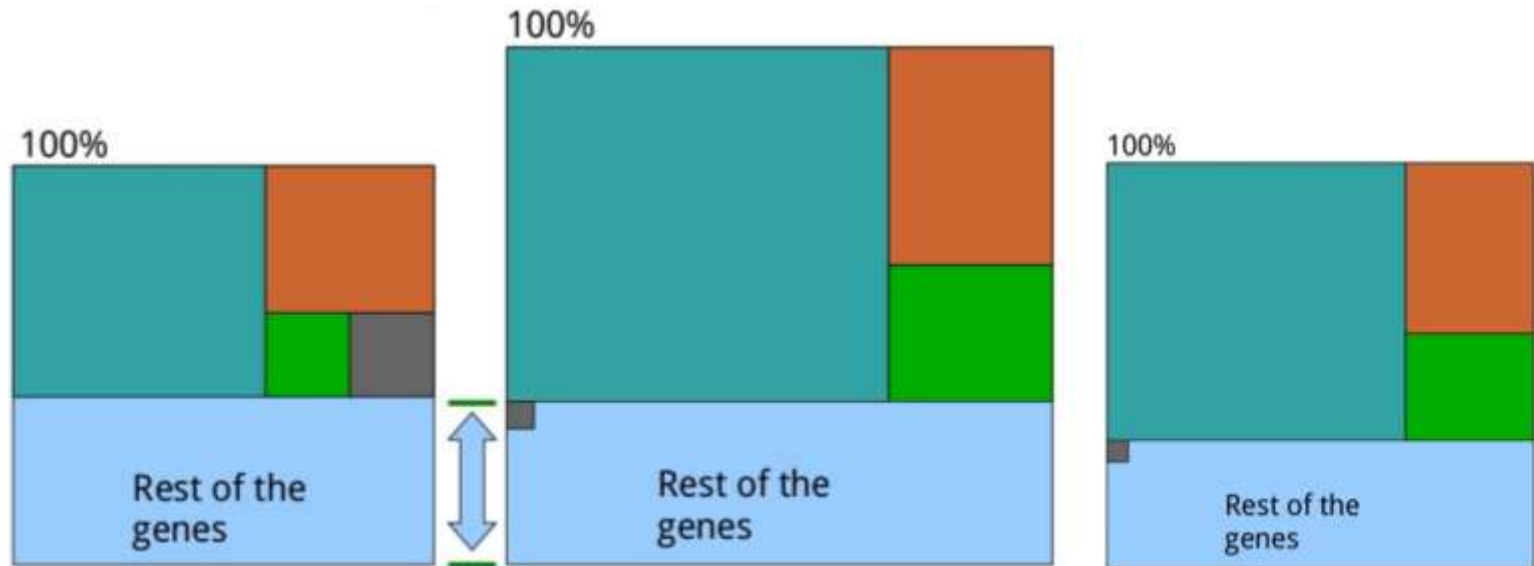


# normalizacja poziomów ekspresji

RPKM:

- Bardzo często wykorzystywane
- Zrozumiałe
- Ale problematyczne:
  1. Zmiany w genach ulegających wysokiej ekspresji powodują zmiany pozostałych genów
  2. Odczyty wielokrotnie mapujące zaburzają zliczenia
  3. Ta sama wartość zmiany RPKM pochodzi od:
    - genu o niskim poziomie i dużej wariacji ekspresji
    - genu o wysokim poziomie i małej wariacji ekspresji

# normalizacja poziomów ekspresji



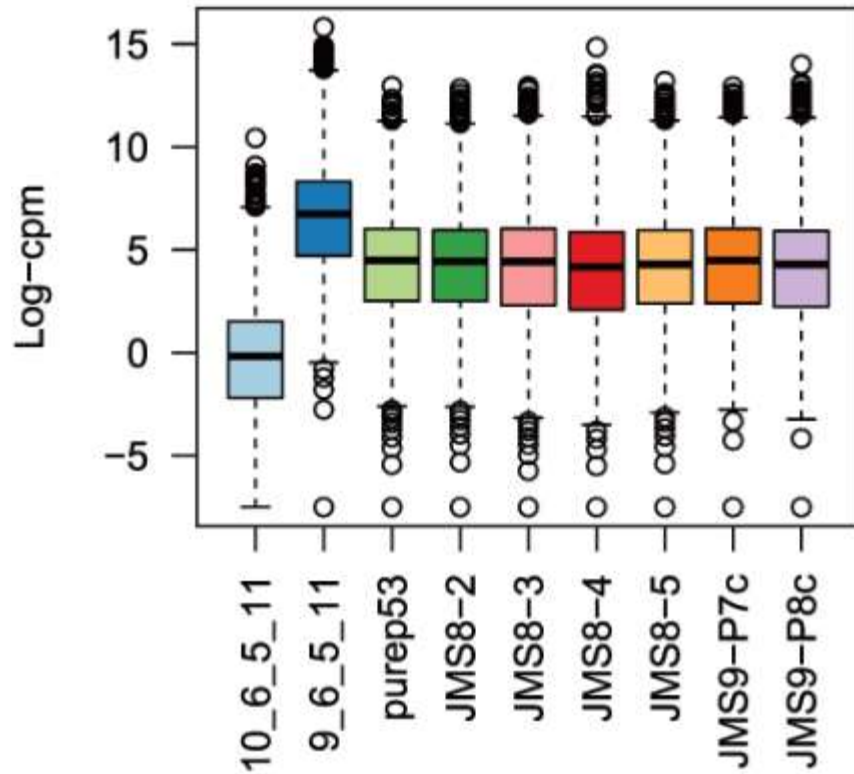


# Metody normalizacji

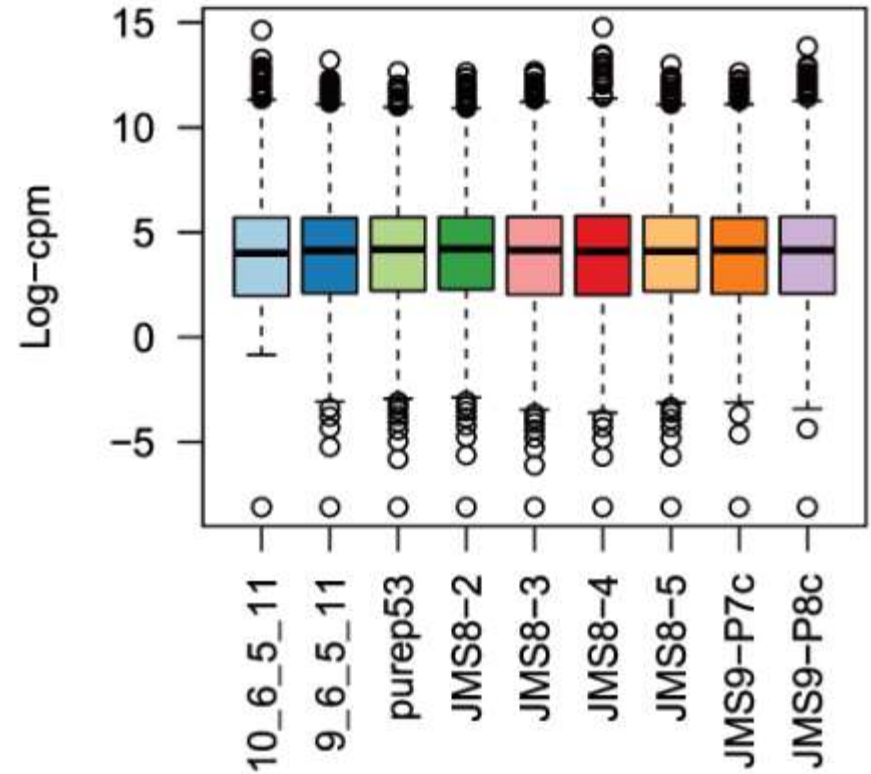
- R/FPKM
  - korekcja różnic w głębokości sekwencjonowania i długości transkryptu
  - cel: porównanie poziomów pomiędzy próbkami i w obrębie próbki
- TMM
  - korekcja różnic w kompozycji próbek spowodowanych pomiarami odstającymi
  - cel: polepszenie precyzji porównań pomiędzy próbkami
- TPM
  - korekcja różnic spowodowanych nieznanymi długościami transkryptów
  - cel: polepszenie precyzji porównań pomiędzy próbkami
- Limma voom (logCPM)
  - stabilizacja wariacji i uniezależnienie jej od średniego poziomu ekspresji

# efekt normalizacije

## A. Example: Unnormalised data

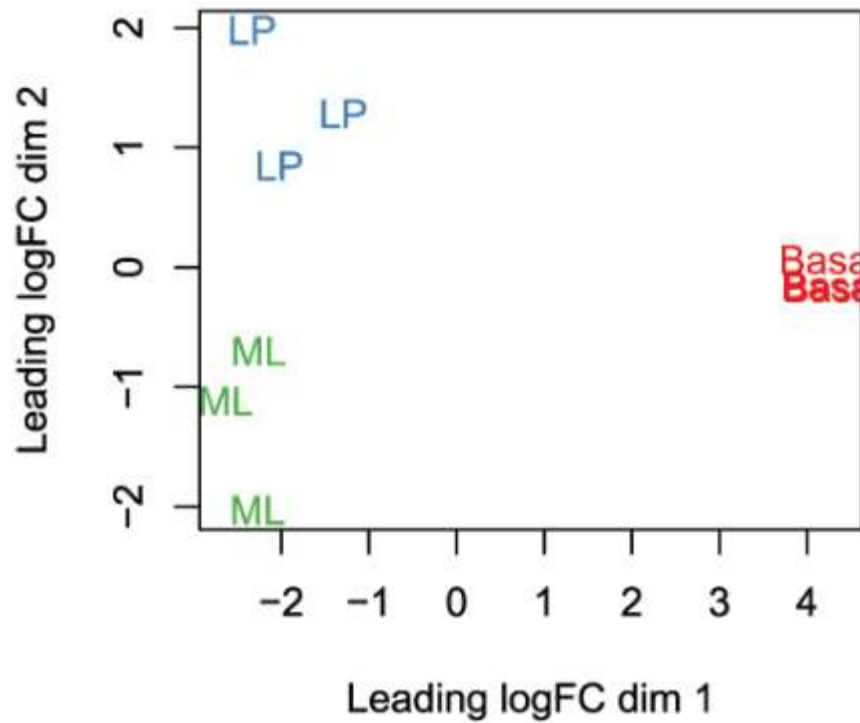


## B. Example: Normalised data

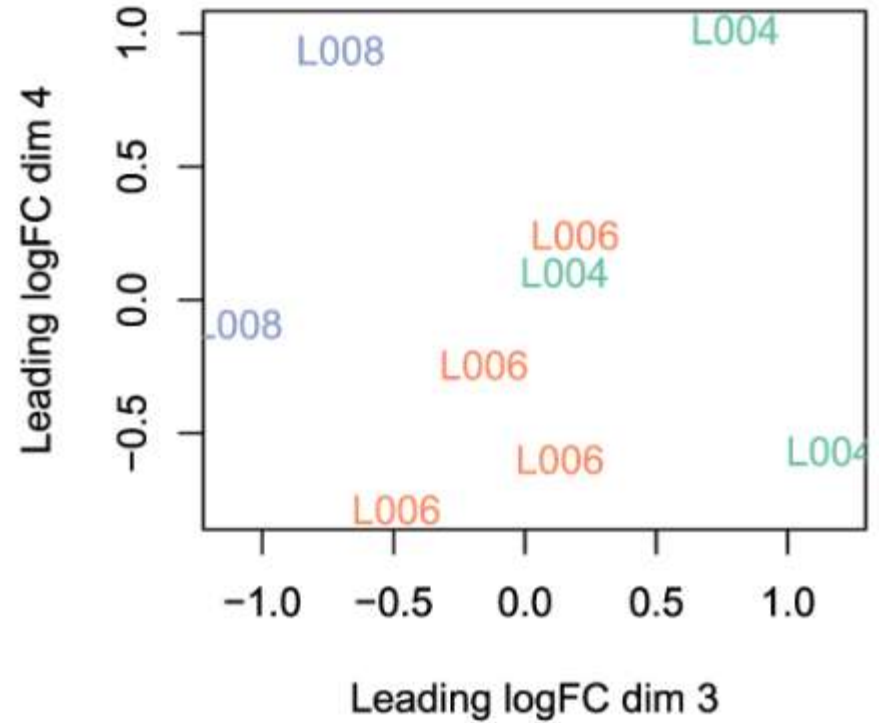


# efekt normalizacije

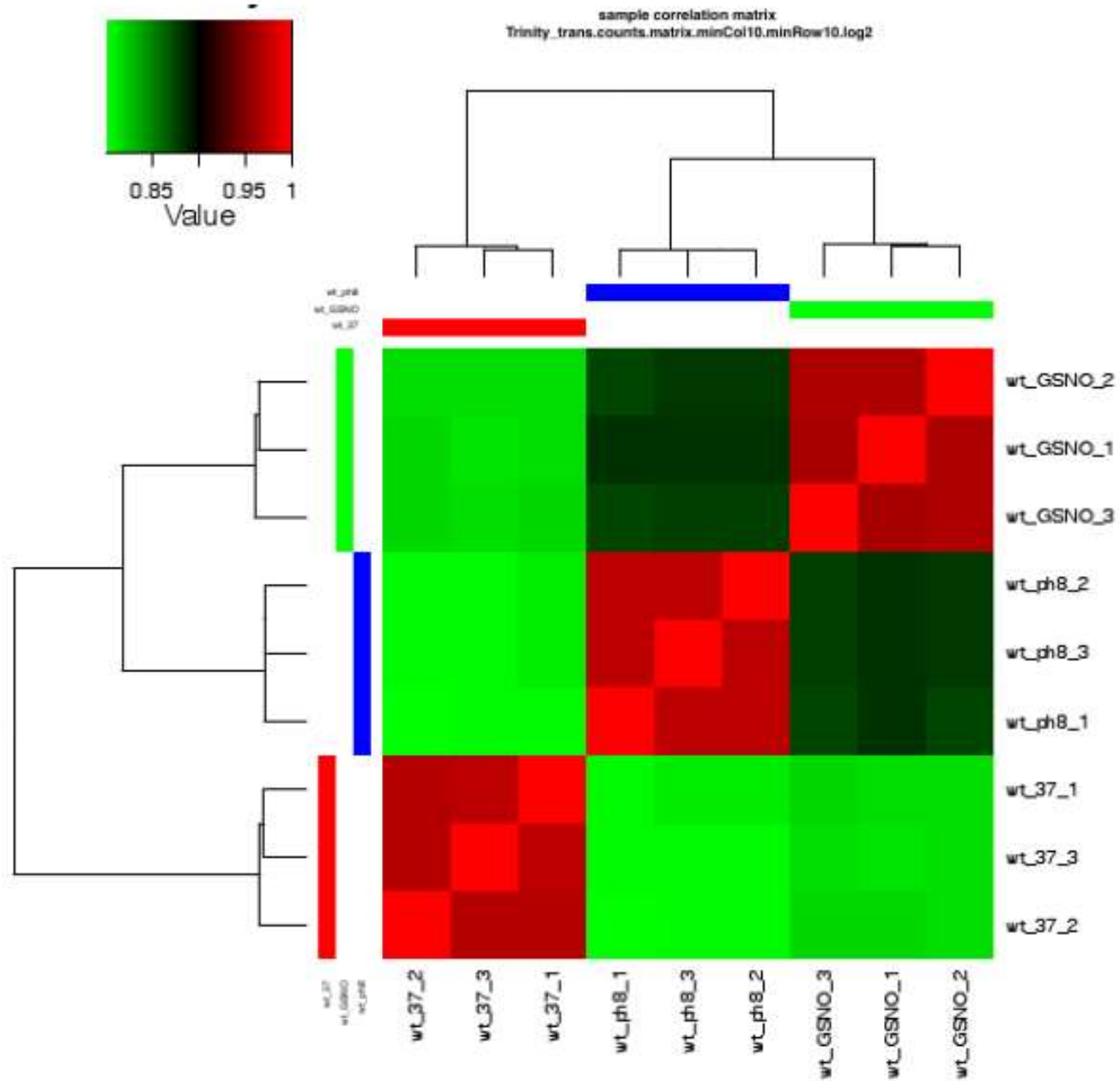
## A. Sample groups



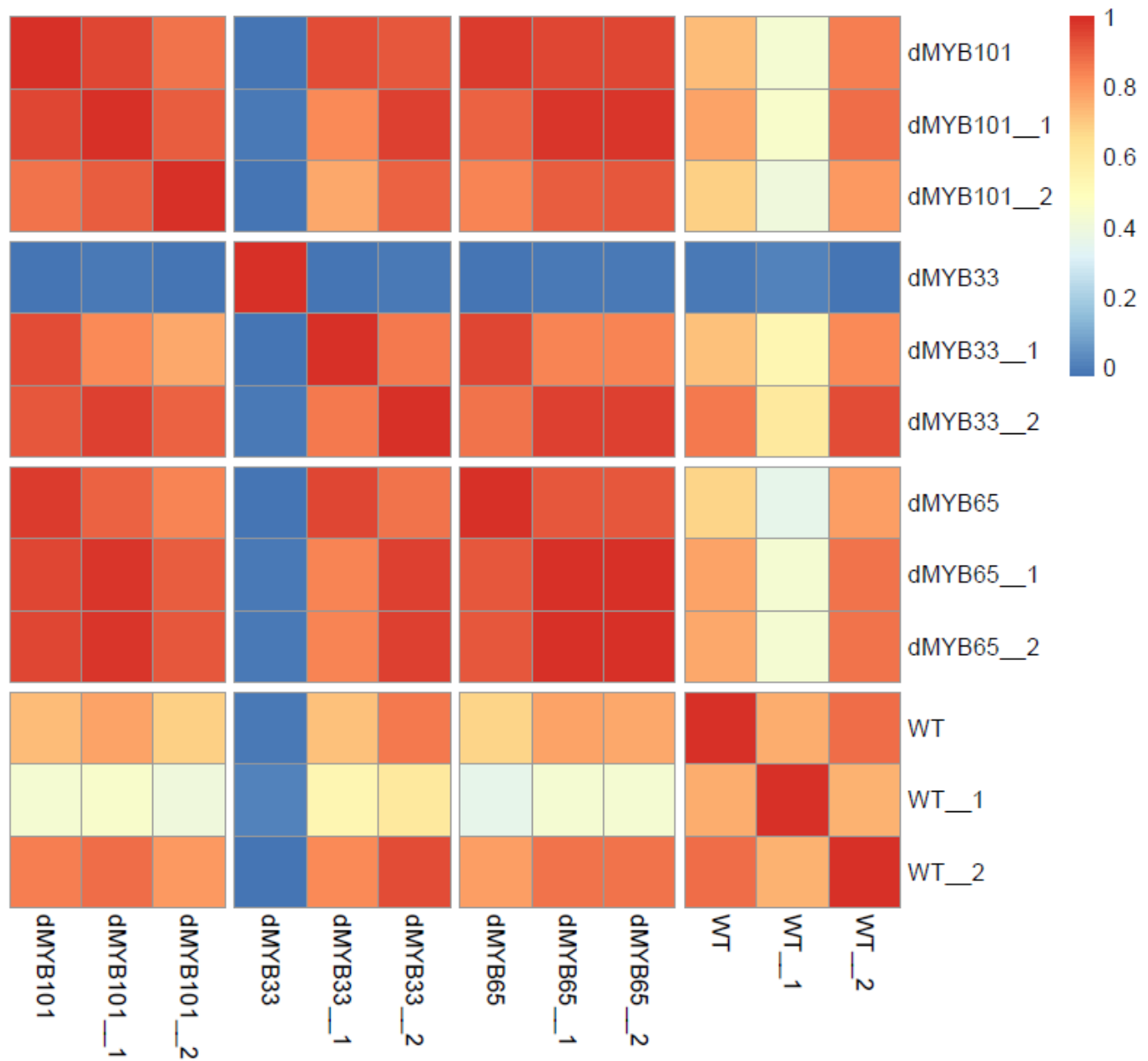
## B. Sequencing lanes



# replikaty



# korelacja próbek - Spearman





# RNA-seq testowanie DE

Metody parametryczne:

- edgeR (negatywna dystrybucja dwumianowa)
- Deseq (negatywna dystrybucja dwumianowa)
- DEGseq (zakłada normalną dystrybucję dla MA)
- baySeq (klasyfikator Bayesa oparty o dystrybucję NB lub Poisson)

Metody nieparametryczne:

- Test Fisher'a
- Cufdiff
- NOlseq (nie potrzebuje replikatów!)

# RNA-seq testowanie DE

Parametry:

## 1. Współczynniki wielkości biblioteki

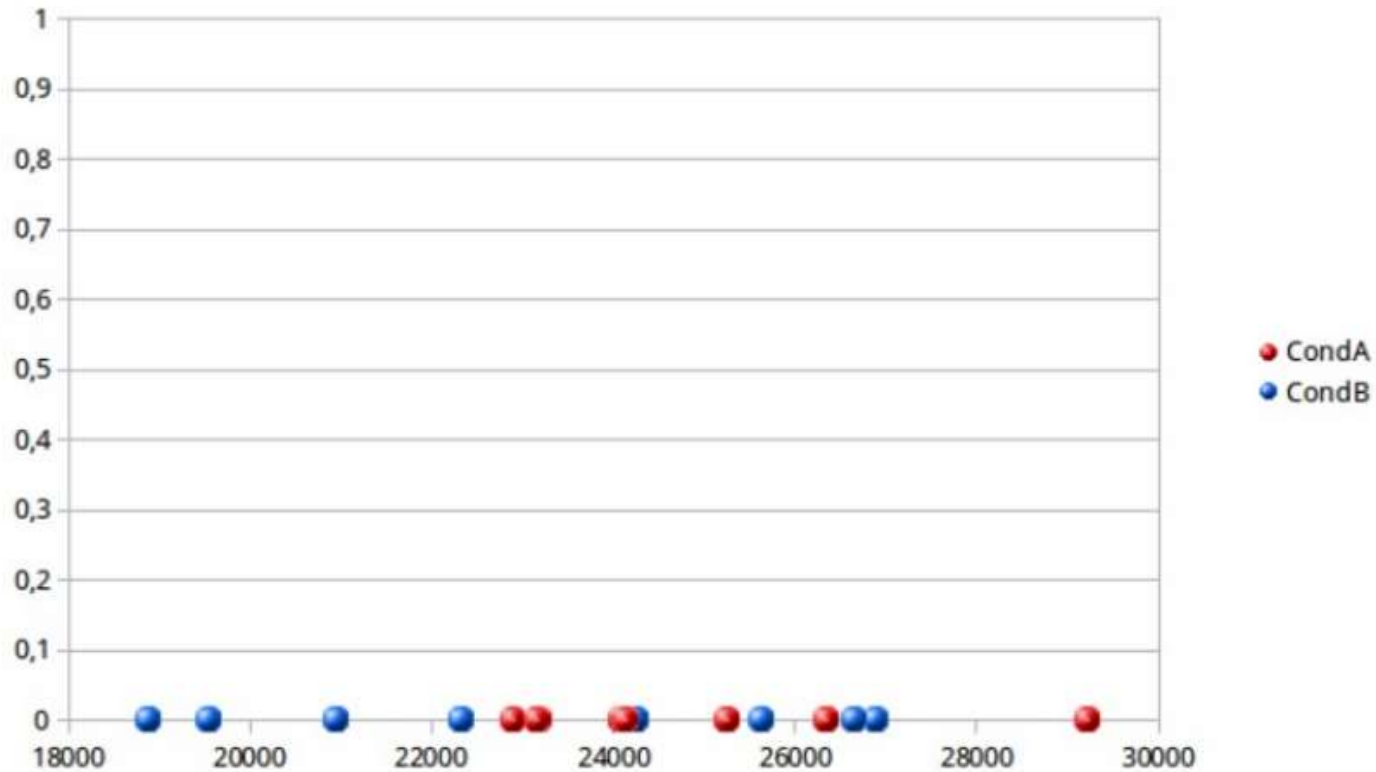
- Obliczane na podstawie obserwowanych różnic w medianie ekspresji
- Brak transformacji odczytów

## 2. Dyspersja w obrębie replikatów

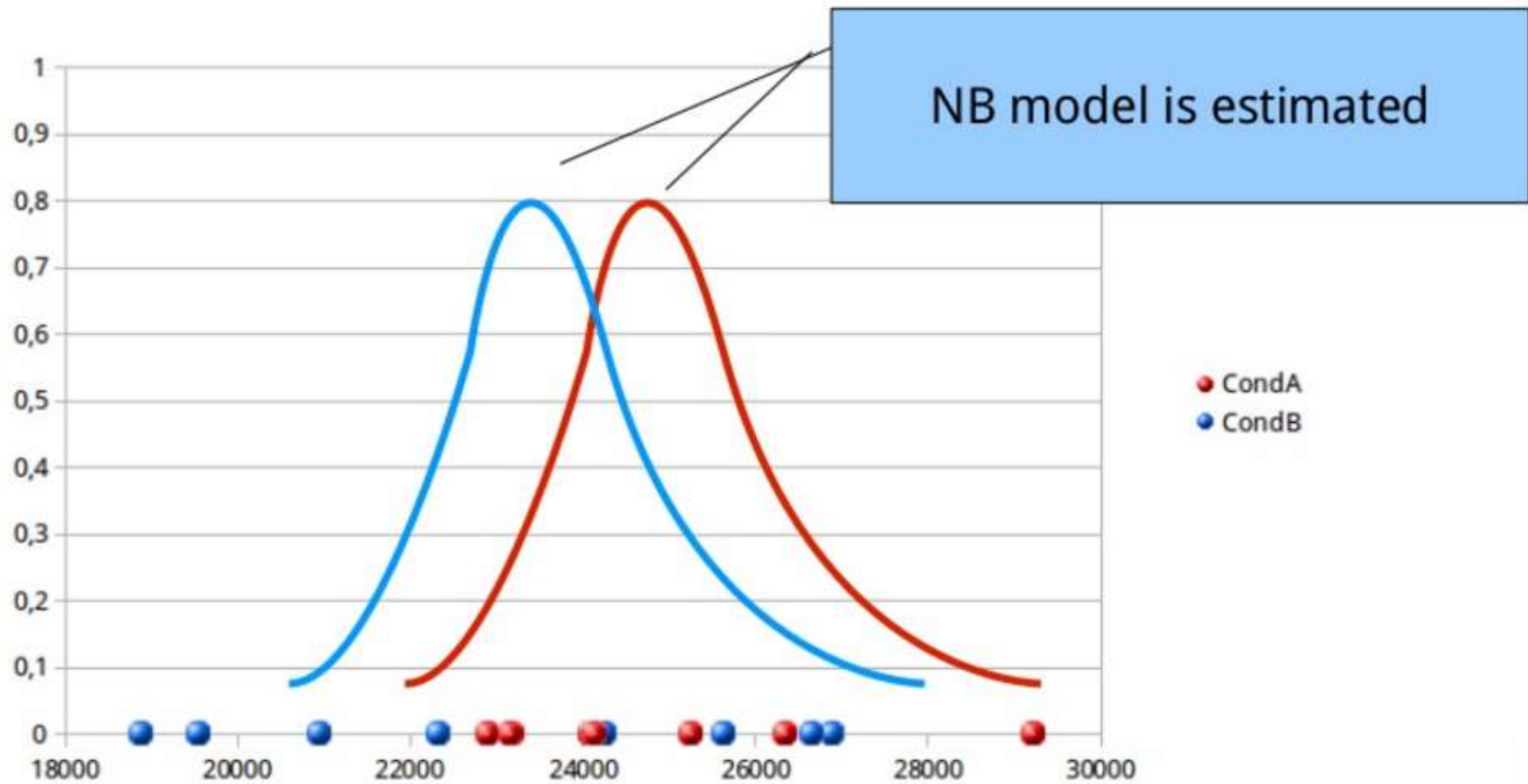
- Liczona dla każdego genu, bądź w binach i uśredniana
- zależna od poziomu ekspresji



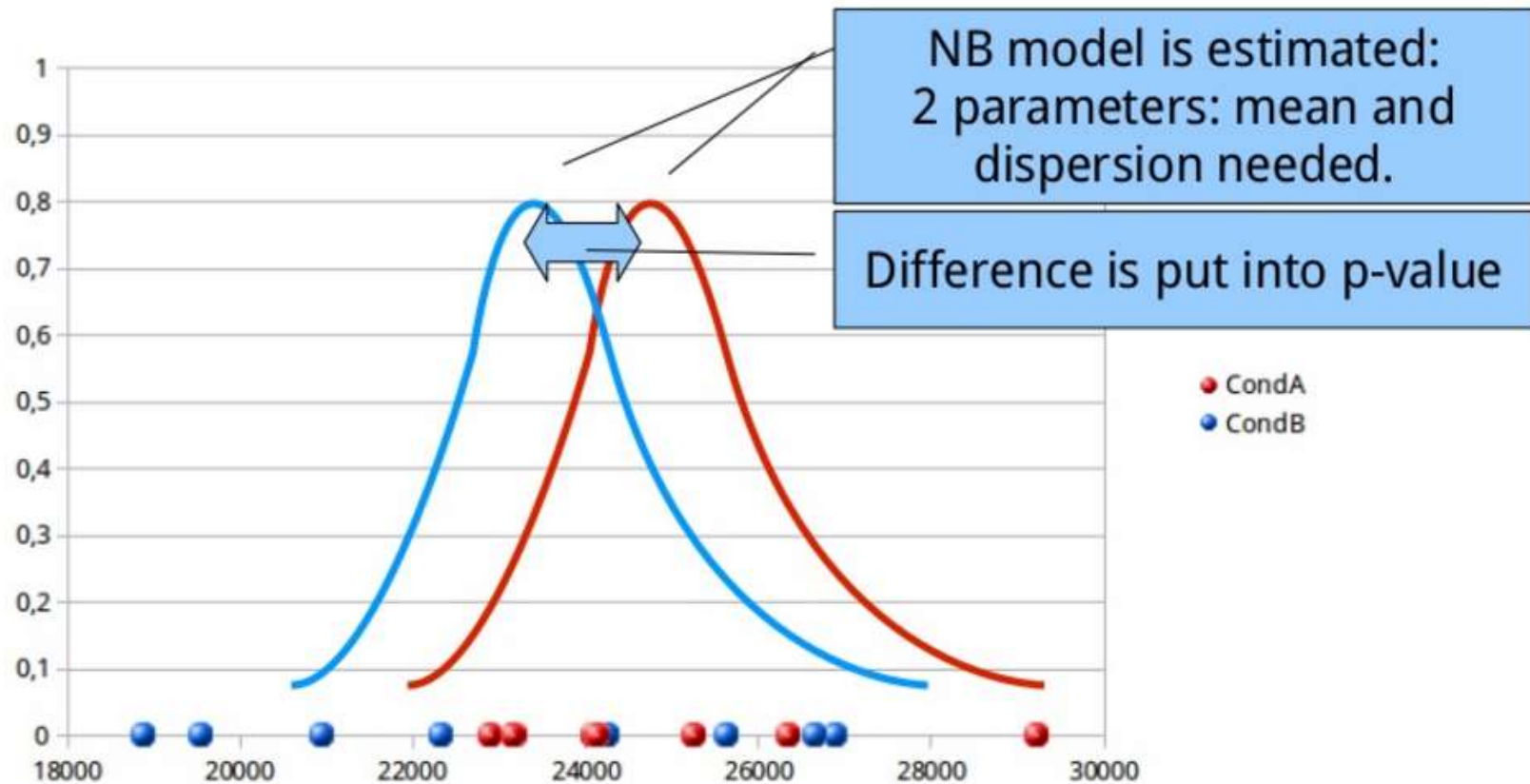
# RNA-seq testowanie DE



# RNA-seq testowanie DE



# RNA-seq testowanie DE



# RNA-seq testowanie DE

