

Analiza danych z wysokoprzepustowego  
sekwencjonowania

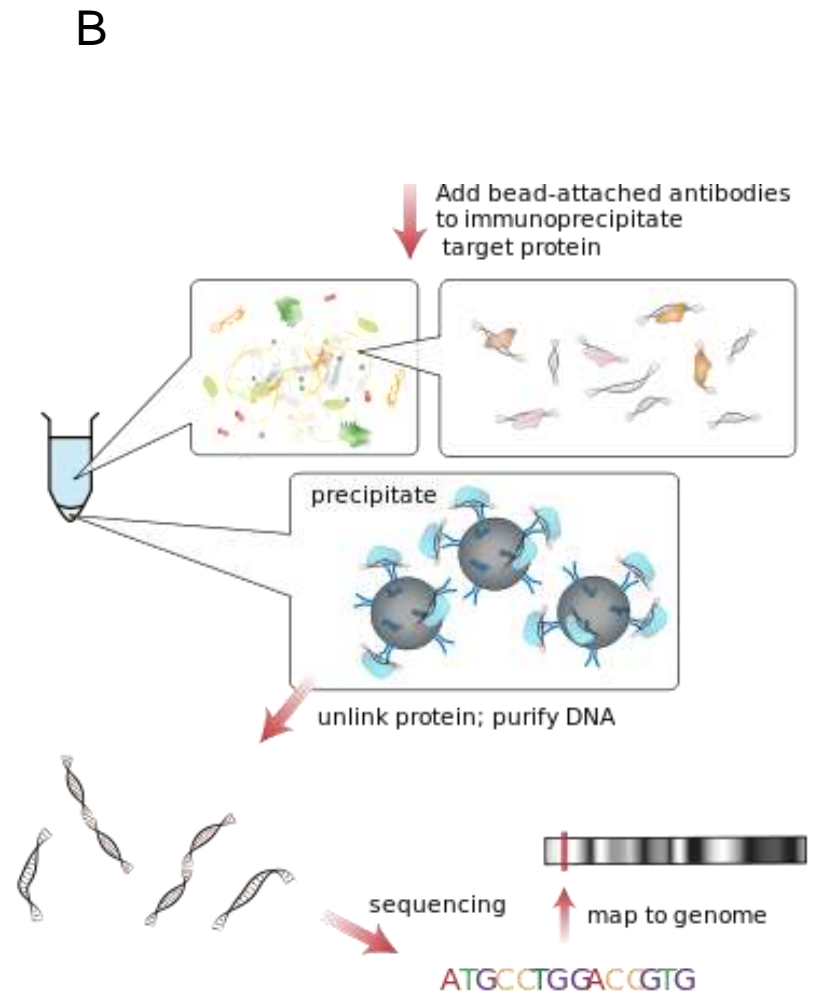
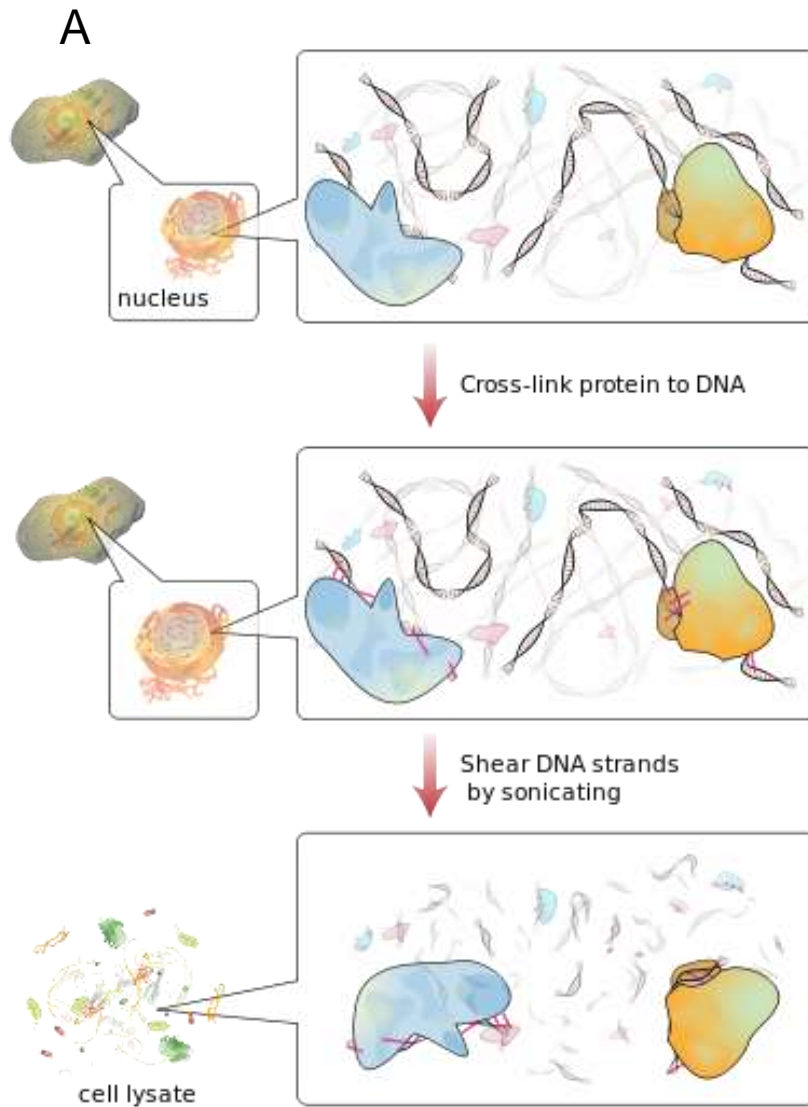
## **W3: Metody oparte o sekwencjonowanie**



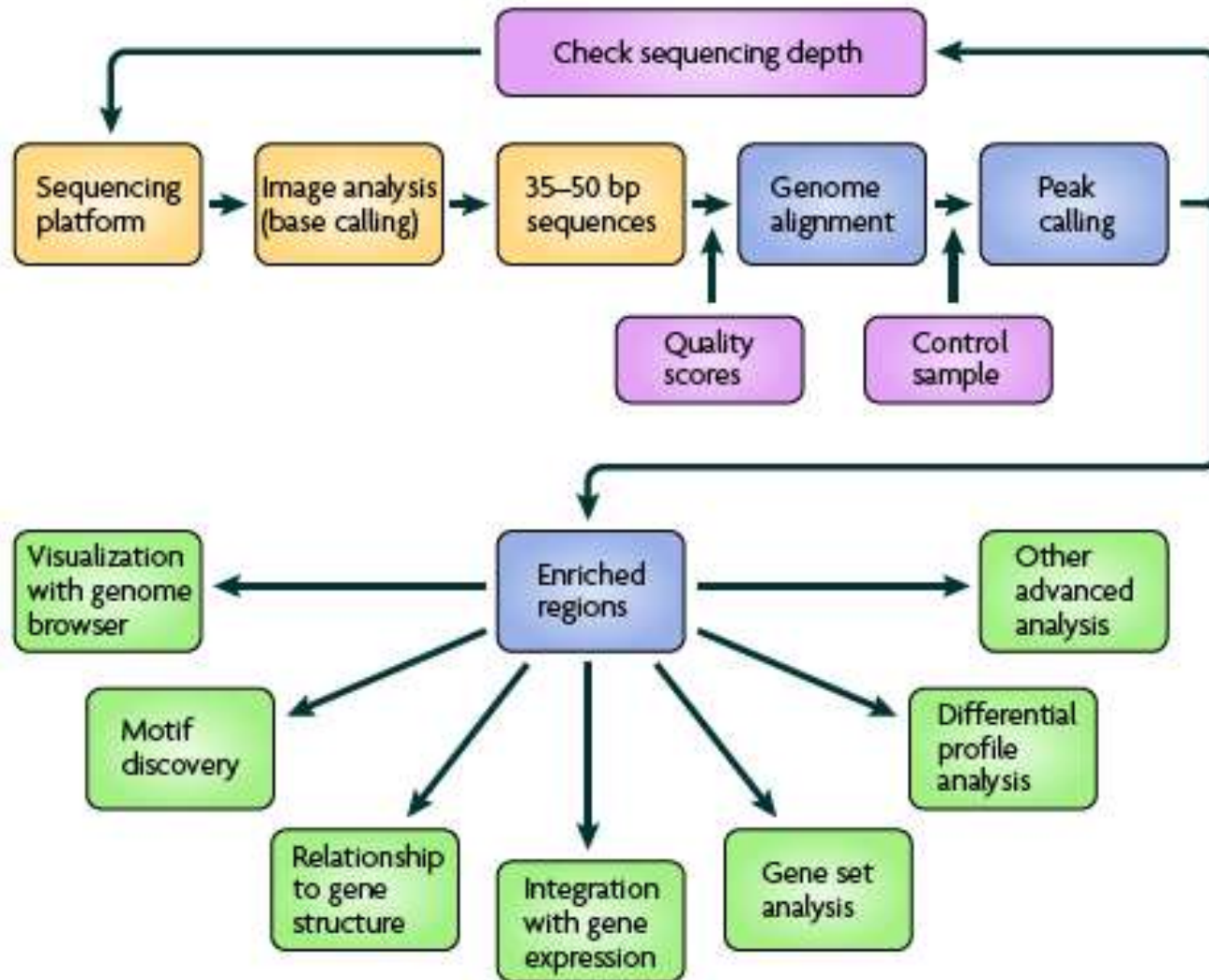
## Chromatin immunoprecipitation and sequencing

- Identyfikacja miejsc wiązania białek do DNA:
  - czynniki transkrypcyjne
  - białka regulatorowe
- Identyfikacja modyfikacji histonów:
  - identyfikacja miejsc aktywnych transkrypcyjnie
- ...

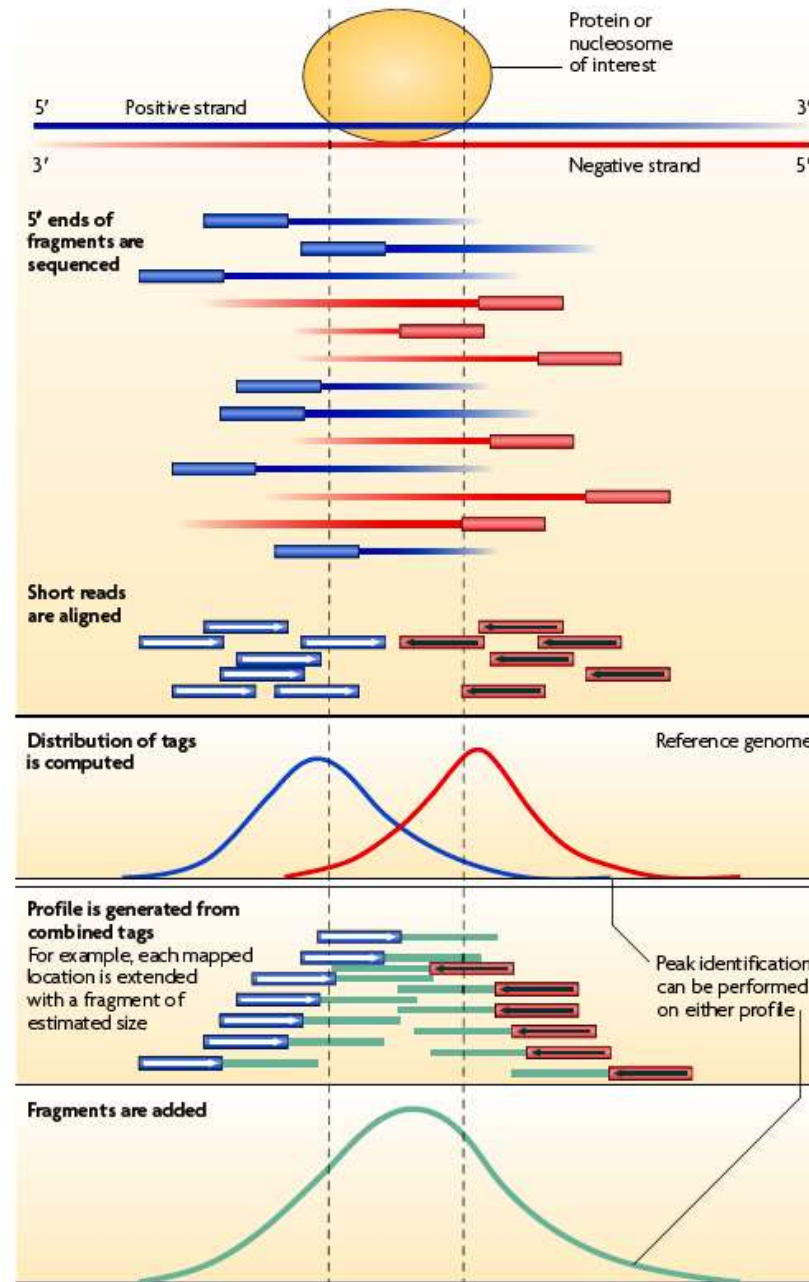
# Chip-Seq



# Chip-Seq

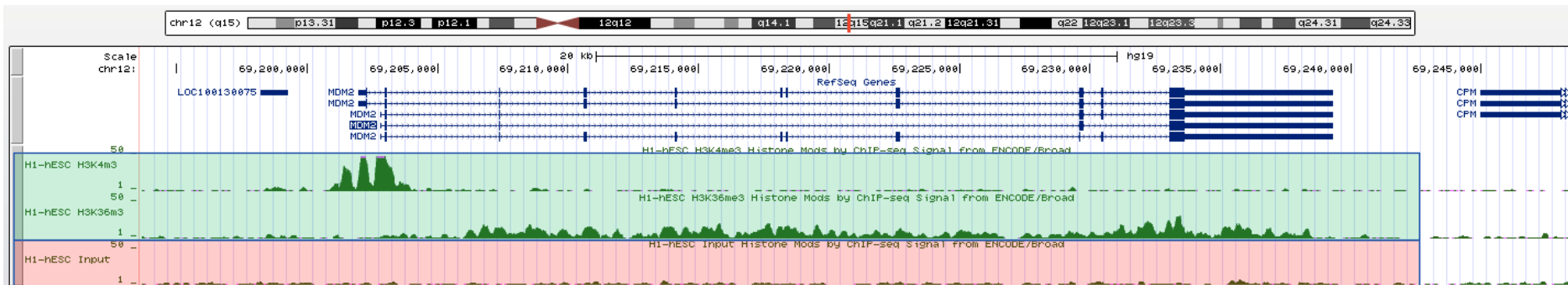


# Chip-Seq



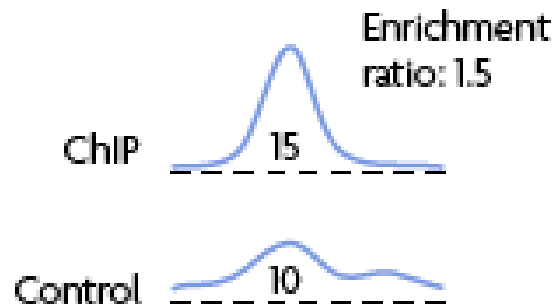
## Identyfikacja rejonów wiązania:

- Poprzez porównanie sygnału uzyskanego dla danych z próbki z przeciwciałem oraz dla kontroli bez przeciwciała
- Szukamy rejonów o statystycznie istotnym wzbogaceniu

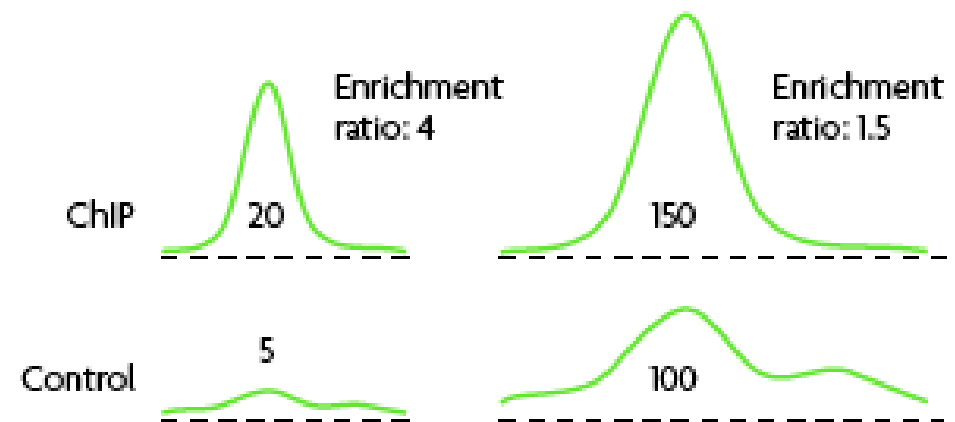


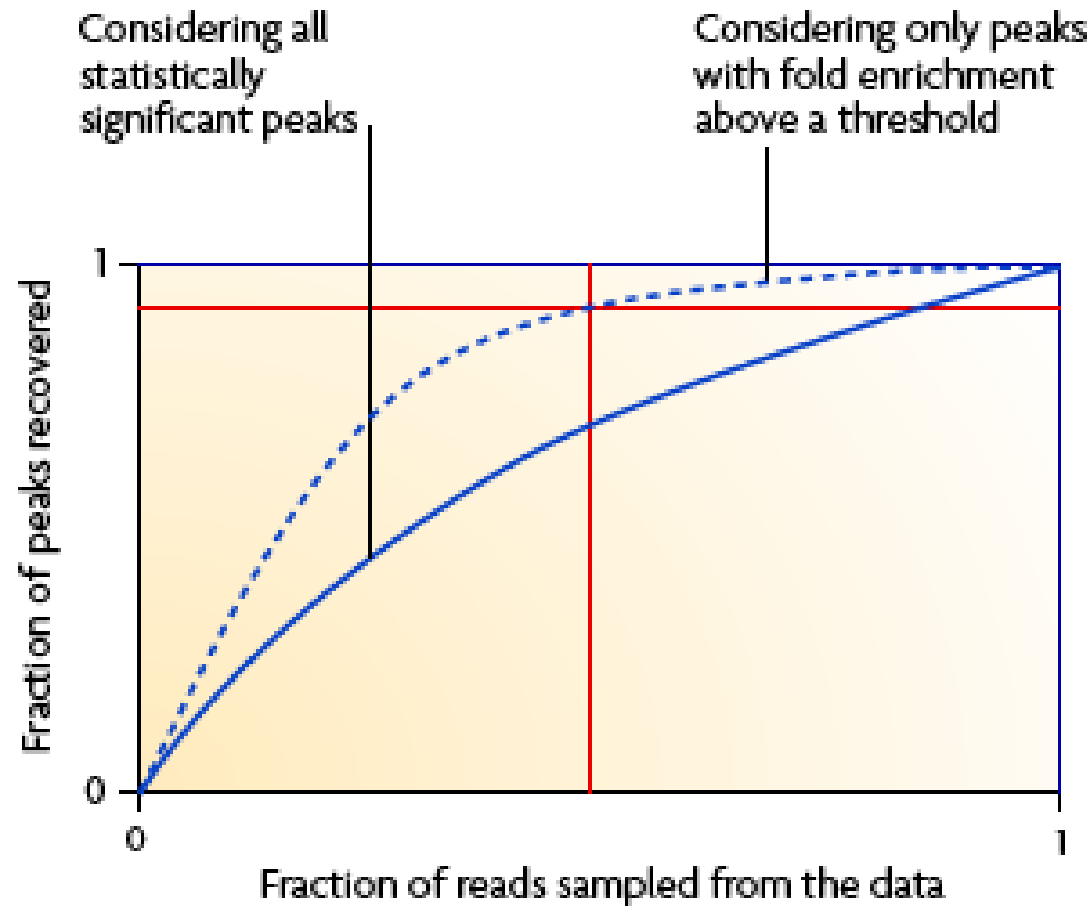
# Chip-Seq

**Ba** Not statistically significant



**Bb** Statistically significant







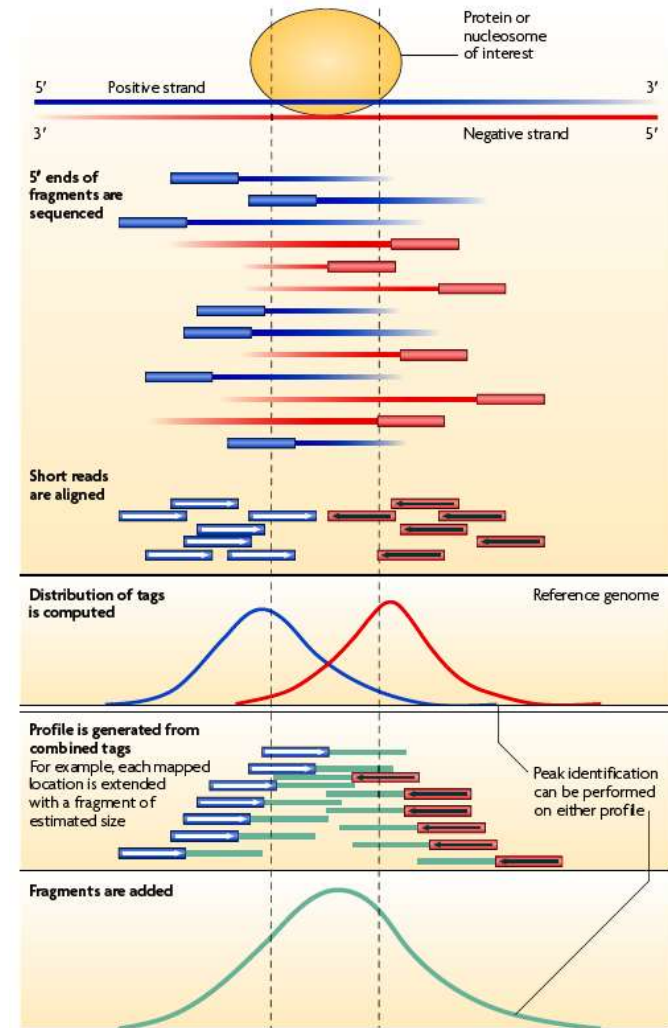
MACS: Model-based Analysis for ChIP-seq

Dwa etapy:

1. Modelowanie przesunięcia sygnału
2. Detekcja wierzchołków

Modelowanie przesunięcia sygnału:

Ponieważ prawdopodobieństwo zsekwencjonowania znaczników z nici plus i minus jest takie samo, rzeczywiste miejsca wiązania wykazują wzbogacenie dwumianowe, o charakterystycznym rozkładzie: na nici plus przed miejscem wiązania oraz na nici minus za miejscem wiązania



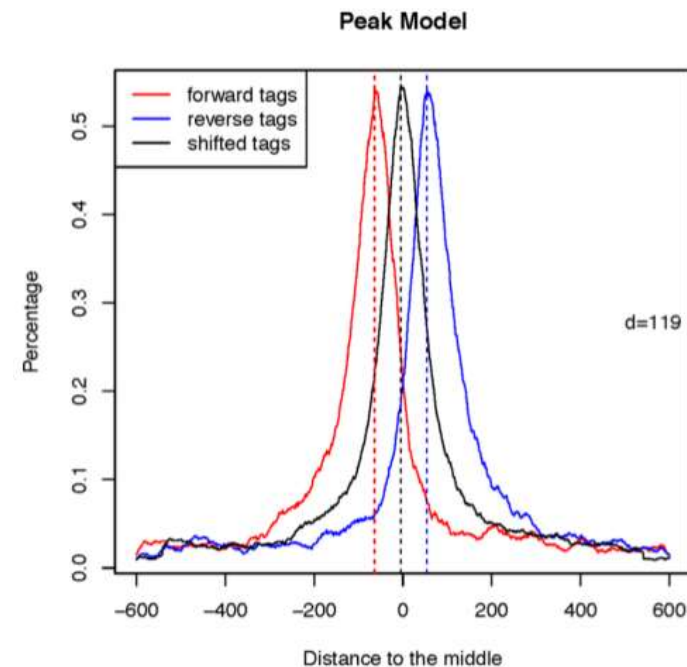
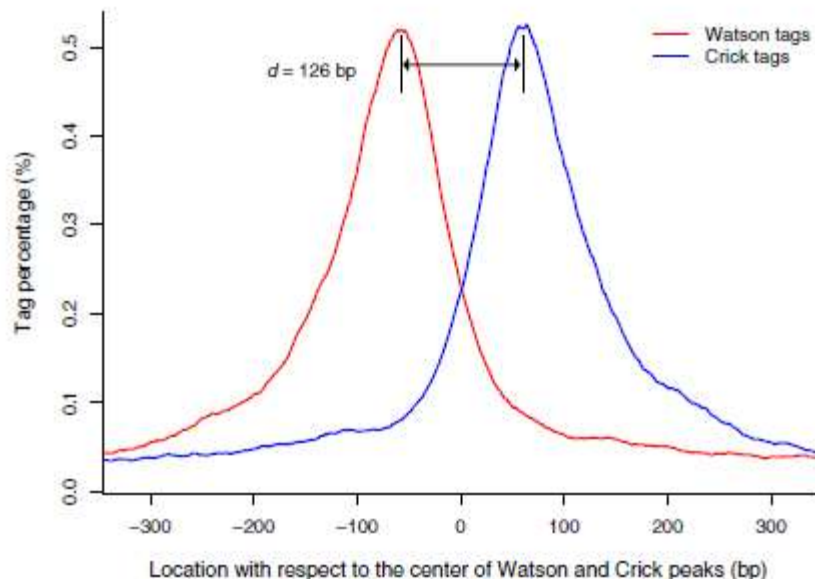
## Modelowanie przesunięcia sygnału:

### Informacje wejściowe:

- wielkość fragmentów DNA użytych do konstrukcji biblioteki
- minimalne wzbogacenie (krotność sygnału w porównaniu z tłem - mfold)

MACS wybiera losowe 1000 rejonów o wysokim wzbogaceniu i na ich podstawie oznacza przesunięcie sygnału:

- $d$  – odległość pomiędzy sygnałami
- $d/2$  – wartość przesunięcia sygnału



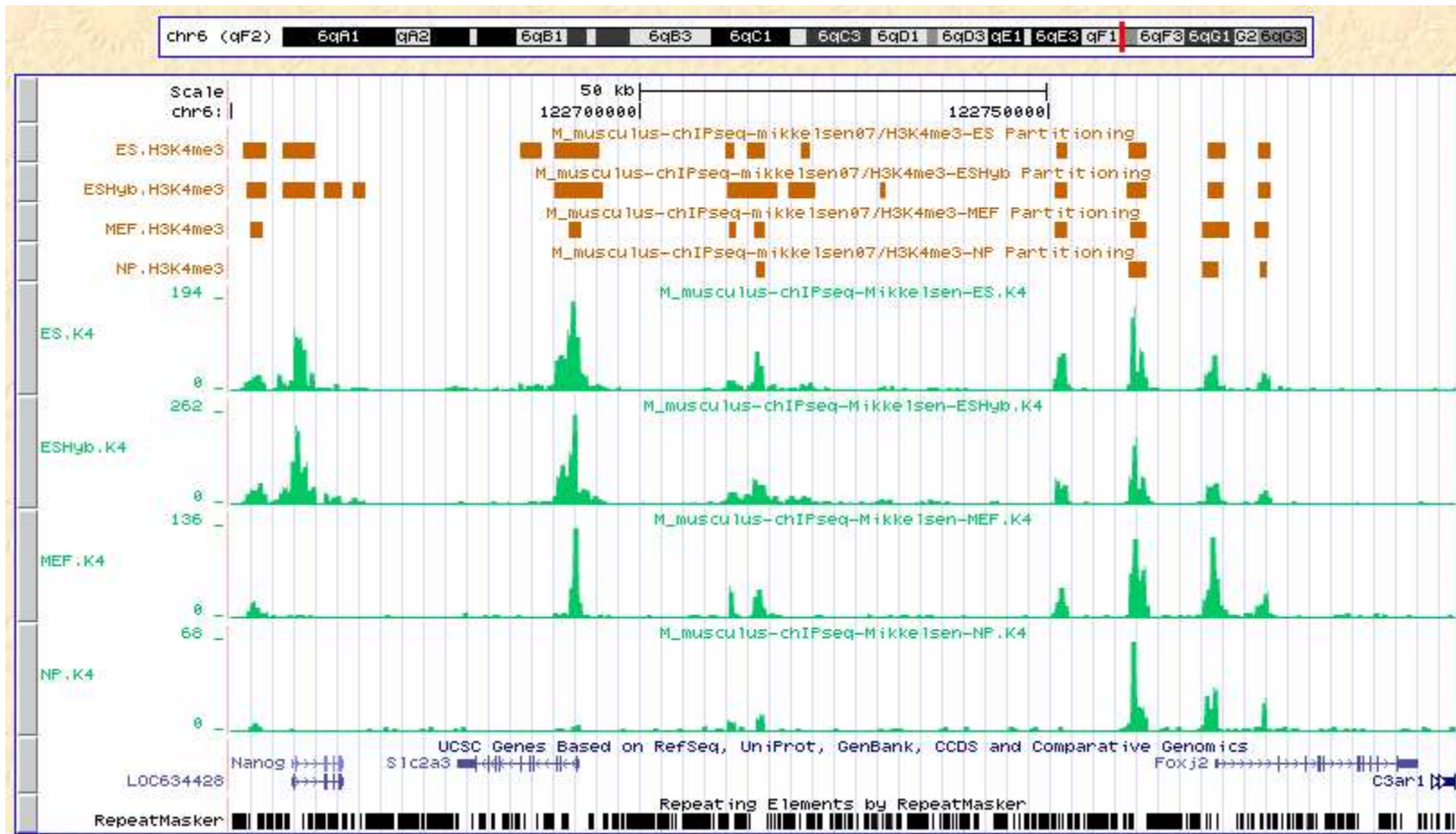
Detekcja miejsc wiązania:

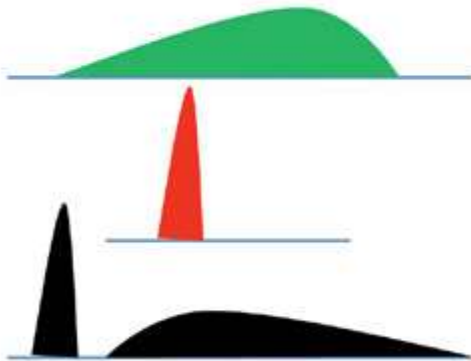
Odczyty zduplikowane są usuwane

Dystrybucja odczytów na genomie może być modelowana za pomocą rozkładu Poisson

Po przesunięciu sygnału o  $d/2$  MACS skanuje genom z użyciem okna o wielkości  $2d$  i wyszukuje rejony o statystycznie znaczącym wzbogaceniu względem kontroli (wartość  $p > 0.05$ )

# Chip-Seq





## Peak calling (choose the right tool)

Type of peak	Example	Representative tools
Broad	H3K27me3	CCAT, SICER
Sharp	CTCF	MACS
Sharp & broad	Pol II	ZINBA

# Chip-Seq


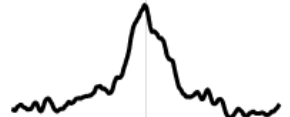
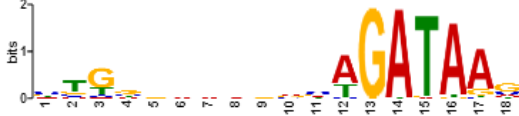
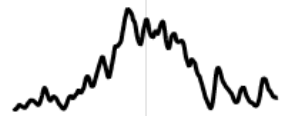



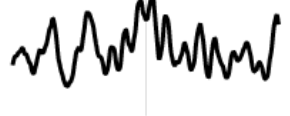
chr15:83139085-83139185	-	23	8.93e-06	TGGAGGCCAG <b>AGGGTGTGGTG</b> GATAAAACCC
chr15:91082416-91082516	-	46	8.93e-06	CTCTGCTGGG <b>AGGGTGTGGTG</b> CCTTTTATCT
chr7:139915433-139915533	+	48	8.93e-06	GAGCGCTGTC <b>AGGGCGGGGCT</b> GTTTTGTAACT
chr16:49777000-49777100	-	89	8.93e-06	G <b>AGGGTGTGGTG</b> GCAACTATTA
chr12:70394726-70394826	-	55	8.93e-06	CAGATAGGAG <b>TGGGTGGGGCA</b> GAGGAAAGTT
chr7:103853792-103853892	+	88	8.93e-06	AATGACAGTA <b>TGGGCGTGGCA</b> AG
chr8:109927244-109927344	+	55	9.76e-06	CGCTAGGGGA <b>AGGGTGGGGCA</b> GCAAAACGGA
chr8:107823858-107823958	+	31	9.76e-06	TGATAAGTTC <b>AGGGTGGGGCA</b> GGTAGGGGAG
chr1:135693776-135693876	-	5	9.76e-06	TGAAGAGGAT <b>AGGGTGGGGCA</b> GGTA
chr1:169322898-169322998	-	60	9.76e-06	CCAGAACCAC <b>AGGGTGGGGCA</b> CCTGCAGGGT
chr5:65203207-65203307	+	42	1.02e-05	TCCGAAGTAC <b>TGGGTGTGGGC</b> AGAATCTTAT
chr3:79383670-79383770	-	50	1.02e-05	CACATTTATT <b>TGGGTGTGGGC</b> CTTCAGCACG
chr11:68709695-68709795	-	50	1.02e-05	CAGTGTCTGT <b>TGGGTGTGGAC</b> CTAGGTCTAC
chr5:96888622-96888722	+	28	1.17e-05	TGTGATGGCA <b>TGGGTGTGGAT</b> GGAGGAGGGA
chr11:104842983-104843083	+	83	1.17e-05	NNNTGCTCTC <b>AGGGTGTGGAC</b> ATTTAGT
chr16:14166506-14166606	+	28	1.17e-05	AGAGGGTAAA <b>AGGGTGTGGGC</b> CTGNNNNNNN
chr17:24289205-24289305	-	69	1.17e-05	GCCTGTGGTA <b>TGGGTGTGGAT</b> TGGGGGACAG
chr7:149571622-149571722	-	38	1.41e-05	CCACCTGACT <b>GGGGTGTGGTT</b> GGTGAGGGAA
chr19:6236070-6236170	+	42	1.41e-05	CCTGCTGGGA <b>GGGGTGTGGTT</b> CACTAGAGGC
chr14:56154883-56154983	-	41	1.41e-05	ACCAGGGGGC <b>AGGGTGTGACC</b> TGCAATGTCC
chr7:99891353-99891453	-	21	1.41e-05	ATTCCTAAAG <b>GGGGTGTGGTT</b> GACCCAAAGT
chr11:53839888-53839988	-	60	1.41e-05	GTAGACAGCT <b>AGGGTGTGGGT</b> CCCACCTGCA
chr7:148354574-148354674	+	44	1.41e-05	TGCCACGAAA <b>TGGGTGTGACT</b> GTGATGAGAG



# Chip-Seq

## MOTIFS

The motifs found by the programs MEME, DREME and CentriMo; clustered by similarity and ordered by E-value.

Motif Found	Discovery/Enrichment Program ?	E-value ?	Known or Similar Motifs ?	Distribution ?
 <p>bits</p> <p>Reverse Complement ⇄ Show 10 More I?</p>	<a href="#">MEME</a>	8.1e-229	<a href="#">Klf4 (MA0039.2)</a> <a href="#">AFT2 (MA0270.1)</a> <a href="#">SP1 (MA0079.1)</a>	
 <p>bits</p> <p>Reverse Complement ⇄ Show 12 More I?</p>	<a href="#">CentriMo</a>	1.3e-045	<a href="#">Tal1::Gata1 (MA0140.1)</a>	
 <p>bits</p> <p>Reverse Complement ⇄</p>	<a href="#">DREME</a>	8.2e-003	<a href="#">TAL1::TCF3 (MA0091.1)</a>	
 <p>bits</p> <p>Reverse Complement ⇄</p>	<a href="#">MEME</a>	1.7e-001	<a href="#">SP1 (MA0079.2)</a> <a href="#">Pax4 (MA0068.1)</a> <a href="#">SP1 (MA0079.1)</a>	



## Zalety:

- Identyfikacja rejonów wiązania białek do DNA w skali genomowej
- Możliwość uzyskania informacji o motywie sekwencyjnym odpowiedzialnym za wiązanie białka

## Wady:

- Czułość metody zależy od siły wiązania białka bądź efektywności cross-link
- Rozdzielczość zależy od analizowanego białka (wysoka dla czynników transkrypcyjnych, niska dla histonów)

## **Na co zwracać uwagę:**

- Spodziewany kształt pików zależy od charakterystyki wiązania DNA przez badane białko
- Czy mamy dostępną kontrolę (eksperyment bez przeciwciała)

Dobór programu do rodzaju pików:

# Chip-Seq

Software tool	Version	Availability	Point- source (peaks)	Broad regions (domains)
BayesPeak [88]	1.10.0	<a href="http://bioconductor.org/packages/release/bioc/html/BayesPeak.html">http://bioconductor.org/packages/release/bioc/html/BayesPeak.html</a>	Yes	
BEADS <sup>§</sup> [84]	1.1	<a href="http://beads.sourceforge.net/">http://beads.sourceforge.net/</a>	Yes	Yes
CCAT [91]	3.0	<a href="http://cmb.gis.a-star.edu.sg/ChIPSeq/paperCCAT.htm">http://cmb.gis.a-star.edu.sg/ChIPSeq/paperCCAT.htm</a>		Yes
CisGenome [56]	2.0	<a href="http://www.biostat.jhsph.edu/~hji/cisgenome/">http://www.biostat.jhsph.edu/~hji/cisgenome/</a>	Yes	
CSAR [85]	1.10.0	<a href="http://bioconductor.org/packages/release/bioc/html/CSAR.html">http://bioconductor.org/packages/release/bioc/html/CSAR.html</a>	Yes	
dPeak	0.9.9	<a href="http://www.stat.wisc.edu/~chungdon/dpeak/">http://www.stat.wisc.edu/~chungdon/dpeak/</a>	Yes	
GPS/GEM [67,18]	1.3	<a href="http://cgs.csail.mit.edu/gps/">http://cgs.csail.mit.edu/gps/</a>	Yes	
HPeak [87]	2.1	<a href="http://www.sph.umich.edu/csg/qin/HPeak/">http://www.sph.umich.edu/csg/qin/HPeak/</a>	Yes	
<b>MACS [17]</b>	<b>2.0.10</b>	<b><a href="https://github.com/taoliu/MACS/">https://github.com/taoliu/MACS/</a></b>	<b>Yes</b>	<b>Yes</b>
NarrowPeaks <sup>§</sup>	1.4.0	<a href="http://bioconductor.org/packages/release/bioc/html/NarrowPeaks.html">http://bioconductor.org/packages/release/bioc/html/NarrowPeaks.html</a>	Yes	
PeakAnalyzer/ PeakSplitter <sup>§</sup> [89]	1.4	<a href="http://www.bioinformatics.org/peakanalyzer">http://www.bioinformatics.org/peakanalyzer</a>	Yes	
PeakRanger [93]	1.16	<a href="http://ranger.sourceforge.net/">http://ranger.sourceforge.net/</a>	Yes	Yes
PeakSeq [24]	1.1	<a href="http://info.gersteinlab.org/PeakSeq">http://info.gersteinlab.org/PeakSeq</a>	Yes	
polyaPeak <sup>§</sup>	0.1	<a href="http://web1.sph.emory.edu/users/hwu30/polyaPeak.html">http://web1.sph.emory.edu/users/hwu30/polyaPeak.html</a>	Yes	
RSEG [92]	0.6	<a href="http://smithlab.usc.edu/histone/rseg/">http://smithlab.usc.edu/histone/rseg/</a>		Yes
SICER [90]	1.1	<a href="http://home.gwu.edu/~wpeng/Software.htm">http://home.gwu.edu/~wpeng/Software.htm</a>		Yes
SIPeS [21]	2.0	<a href="http://gmdd.shgmo.org/Computational-Biology/ChIP-Seq/download/SIPeS">http://gmdd.shgmo.org/Computational-Biology/ChIP-Seq/download/SIPeS</a>	Yes	
SISSRs [19]	1.4	<a href="http://sisrs.rajajothi.com/">http://sisrs.rajajothi.com/</a>	Yes	
SPP [9]	1.1	<a href="http://compbio.med.harvard.edu/Supplements/ChIP-seq/">http://compbio.med.harvard.edu/Supplements/ChIP-seq/</a>	Yes	Yes
USeq [97]	8.5.1	<a href="http://sourceforge.net/projects/useq/">http://sourceforge.net/projects/useq/</a>	Yes	
ZINBA [86]	2.02.03	<a href="http://code.google.com/p/zinba/">http://code.google.com/p/zinba/</a>	Yes	Yes

# Chip-Seq – metody normalizacji

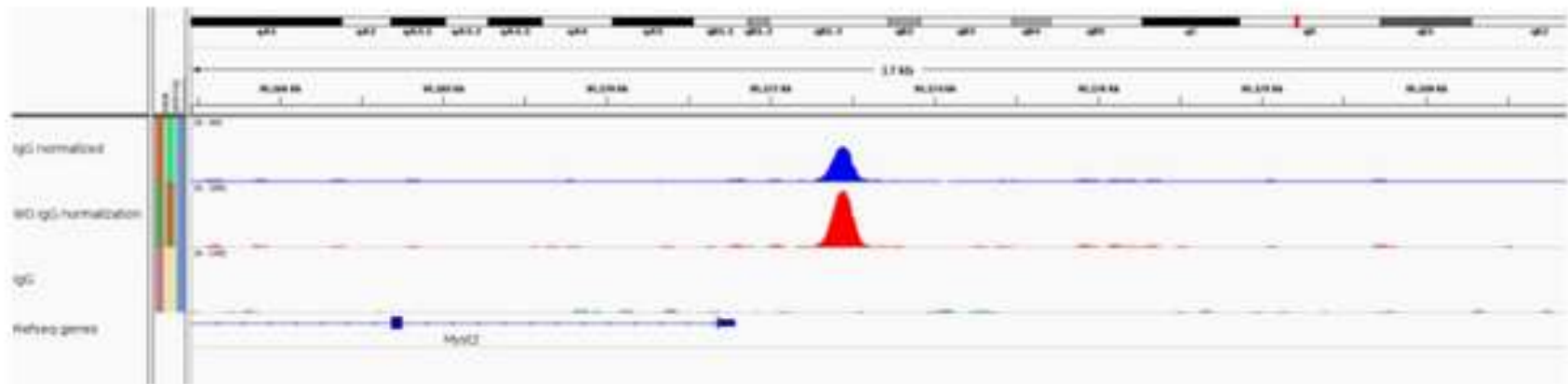
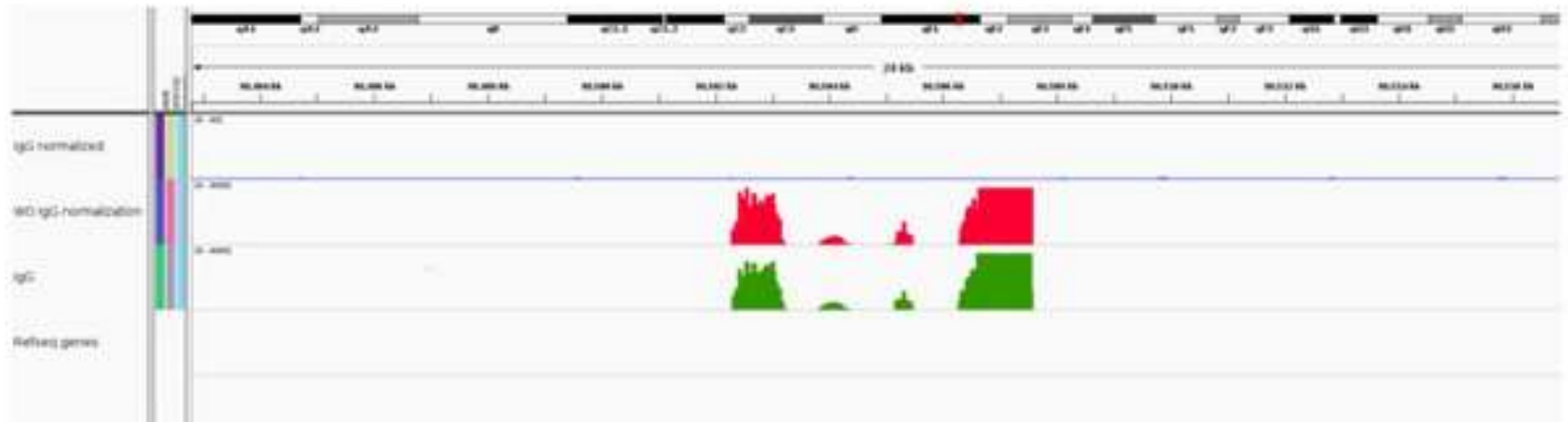
<b>Linear Algorithms</b>	<b>Focus on negative control sample</b>	<b>Publication Time</b>
PeakSeq [24]	Yes	Jan 2009
Cisgenome [56]	Yes	Nov 2008
MACS [17]	Yes	Sept 2008
USeq [97]	Yes	Dec 2008
RPKM [98]	No	May 2008

<b>Non-linear Algorithms</b>	<b>Focus on negative control sample</b>	<b>Publication Time</b>
Locally weighted regression with respect to mean and variance [28]	No	Jan 2009
MAnorm [34]	No	March 2012
POLYPHEMUS [37] (Quantile and locally weighted regression)	No	Feb 2012

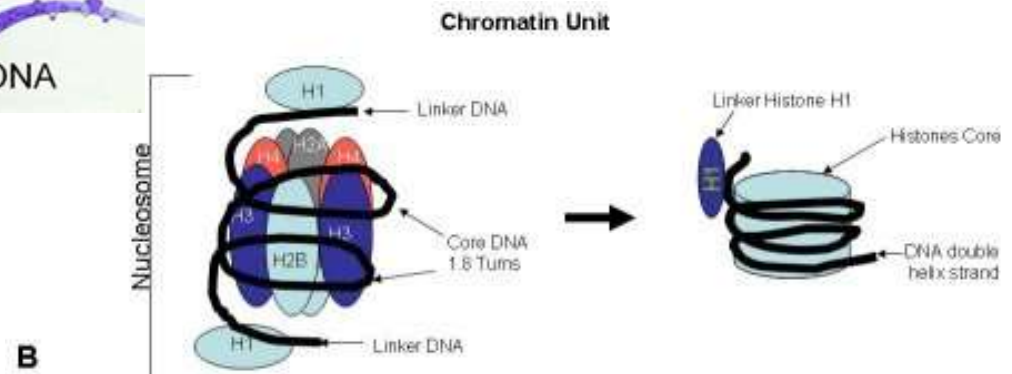
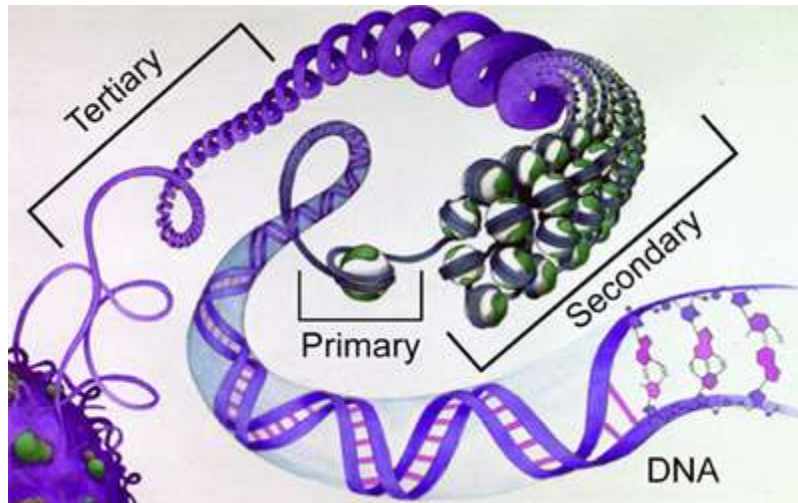
# Chip-Seq – różnicowe wiązanie

Software tool	Availability	Notes
<b>ChIPDiff [36]</b>	<a href="http://cmb.gis.a-star.edu.sg/ChIPSeq/paperChIPDiff.htm">http://cmb.gis.a-star.edu.sg/ChIPSeq/paperChIPDiff.htm</a>	Differential histone modification sites using a hidden Markov model
<b>Comparative ChIP-seq [25]</b>	<a href="http://www.starklab.org/data/bardet_natprotoc_2011/">http://www.starklab.org/data/bardet_natprotoc_2011/</a>	Fold change ratio between normalized peak heights
<b>DBChIP [33]</b>	<a href="http://pages.cs.wisc.edu/~kliang/DBChIP/">http://pages.cs.wisc.edu/~kliang/DBChIP/</a>	Assigns uncertainty measures in a test of non-differential binding (uses edgeR)
<b>DESeq<sup>s</sup> [31]</b>	<a href="http://www.bioconductor.org/packages/release/bioc/html/DESeq.html">http://www.bioconductor.org/packages/release/bioc/html/DESeq.html</a>	Test based on a model using the negative binomial distribution
<b>DiffBind</b>	<a href="http://www.bioconductor.org/packages/release/bioc/html/DiffBind.html">http://www.bioconductor.org/packages/release/bioc/html/DiffBind.html</a>	Differential binding affinity analysis (uses edgeR and DESeq)
<b>DIME [35]</b>	<a href="http://cran.r-project.org/web/packages/DIME/">http://cran.r-project.org/web/packages/DIME/</a>	Differential identification using mixtures ensemble
<b>edgeR<sup>s</sup> [32]</b>	<a href="http://www.bioconductor.org/packages/release/bioc/html/edgeR.html">http://www.bioconductor.org/packages/release/bioc/html/edgeR.html</a>	Empirical Bayes estimation and exact tests based on the negative binomial distribution
<b>MACS [17] (version 2)</b>	<a href="https://github.com/taoliu/MACS/">https://github.com/taoliu/MACS/</a>	Differential peak detection based on paired four bedGraph files
<b>MAnorm [34]</b>	<a href="http://bcb.dfci.harvard.edu/~gcyuan/MAnorm/MAnorm.htm">http://bcb.dfci.harvard.edu/~gcyuan/MAnorm/MAnorm.htm</a>	Robust regression to derive a linear model
<b>MMDiff</b>	<a href="http://bioconductor.org/packages/release/bioc/html/MMDiff.html">http://bioconductor.org/packages/release/bioc/html/MMDiff.html</a>	Differences in shape using Kernel methods
<b>NarrowPeaks</b>	<a href="http://bioconductor.org/packages/release/bioc/html/NarrowPeaks.html">http://bioconductor.org/packages/release/bioc/html/NarrowPeaks.html</a>	Shape-based analysis of variation using functional PCA
<b>POLYPHEMUS [37]</b>	<a href="http://cran.r-project.org/web/packages/polyphemus/">http://cran.r-project.org/web/packages/polyphemus/</a>	Non-linear normalization on RNA Pol II profiling

# Chip-Seq – wpływ kontroli

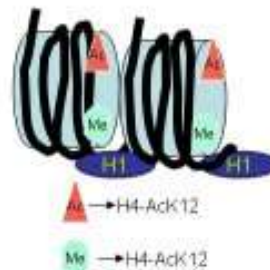


# Struktura DNA

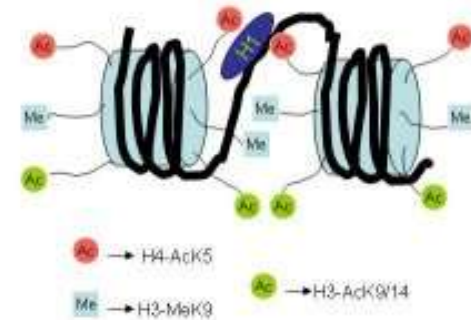


**B**

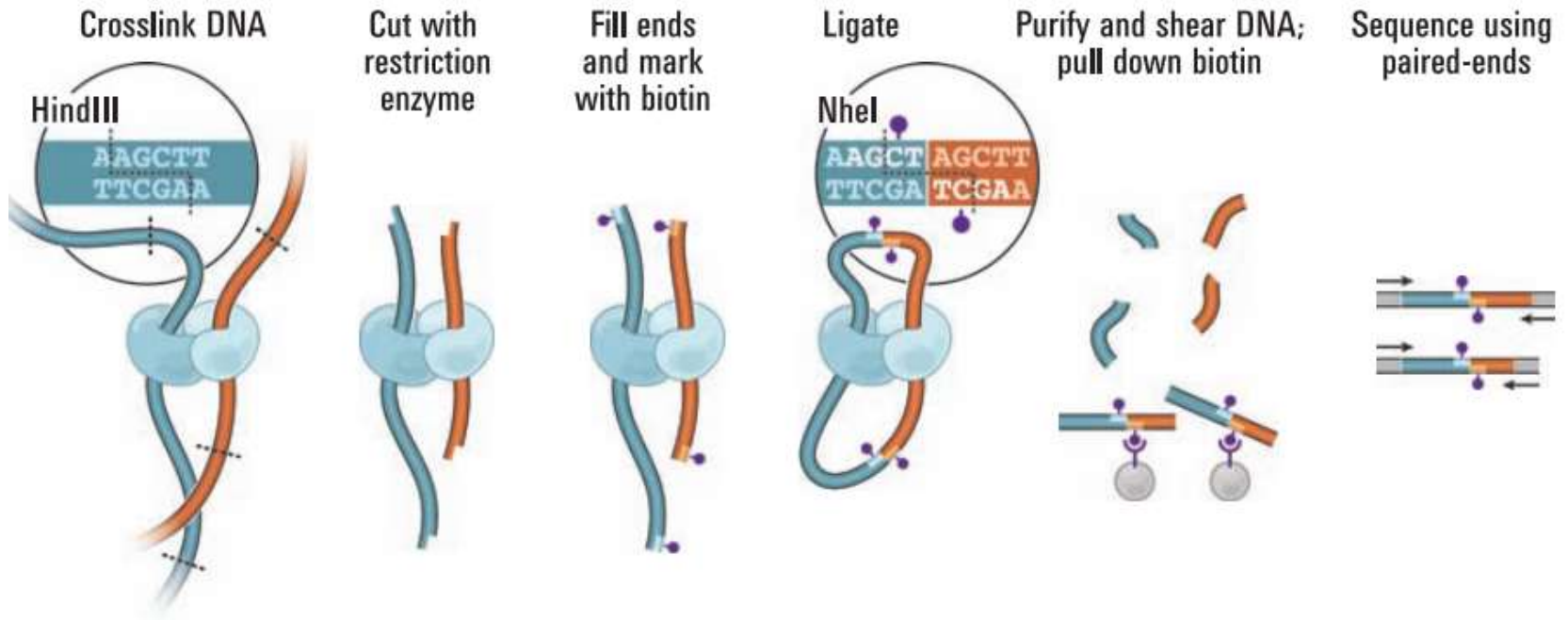
**Closed/Inactive Chromatin**



**Open/Active Chromatin**

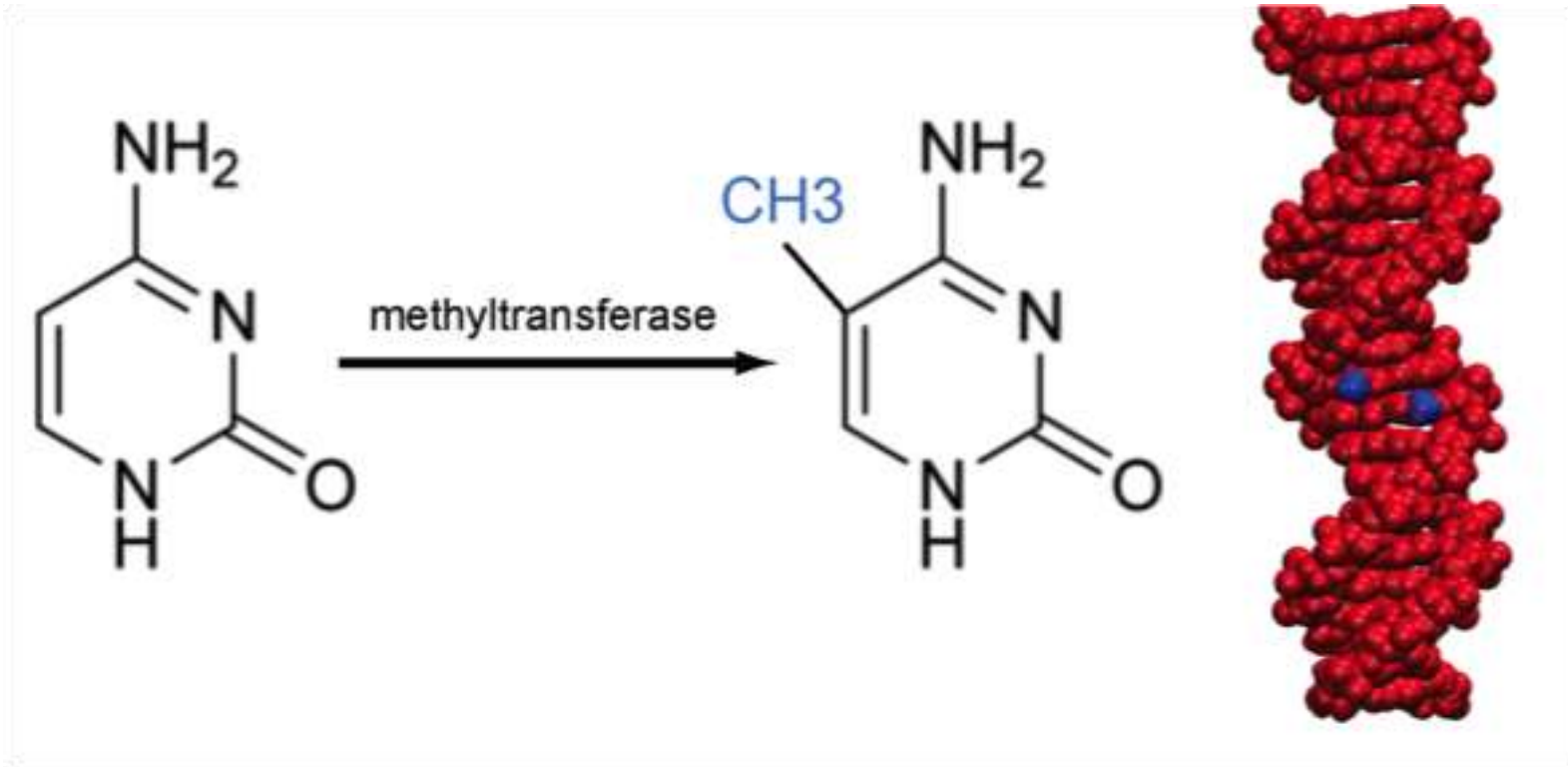


# Chromatin interactions mapping (Hi-C)

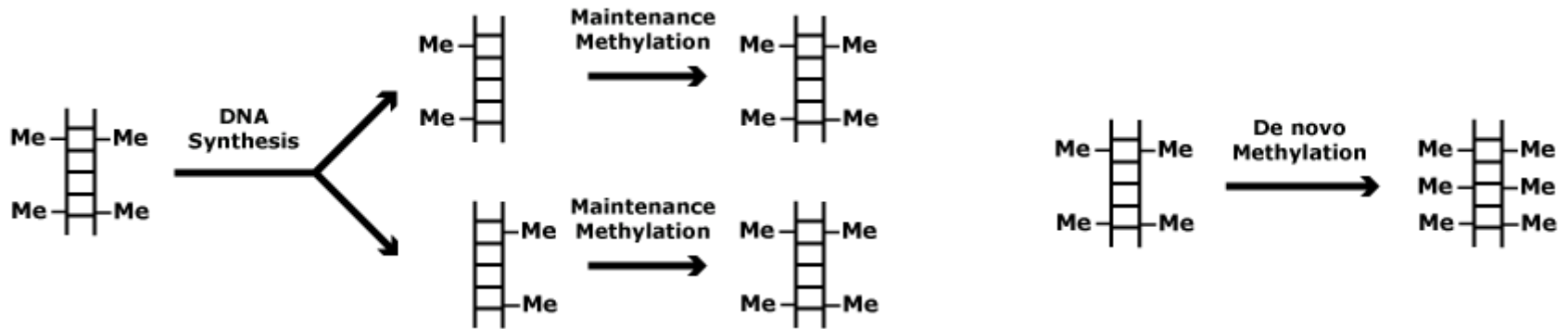
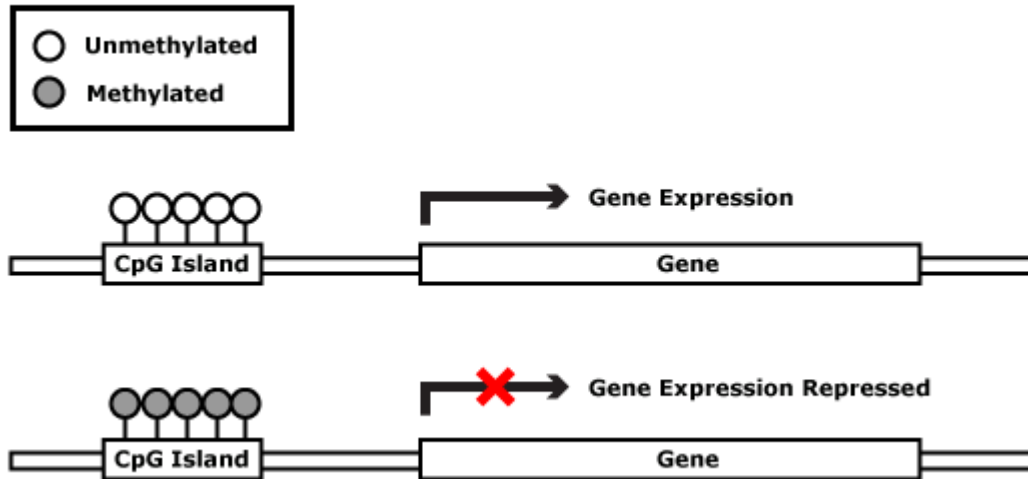




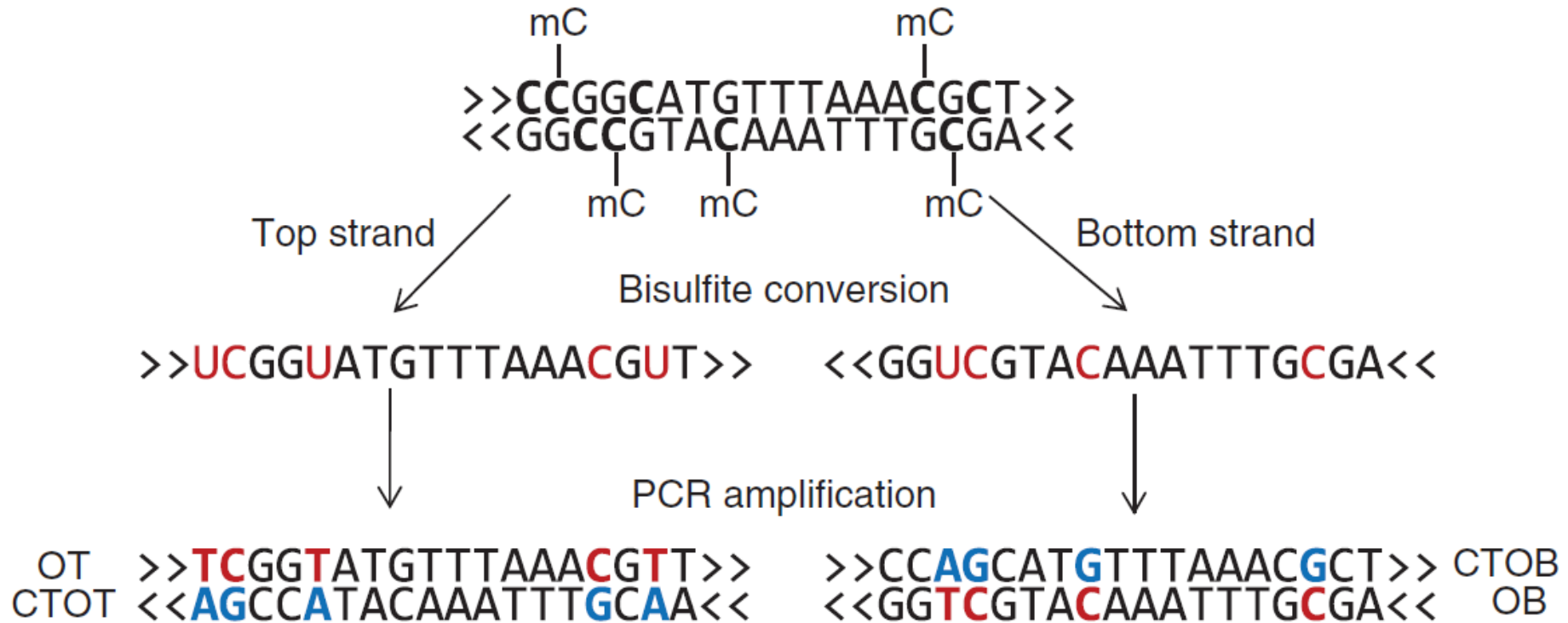
# metylacja DNA



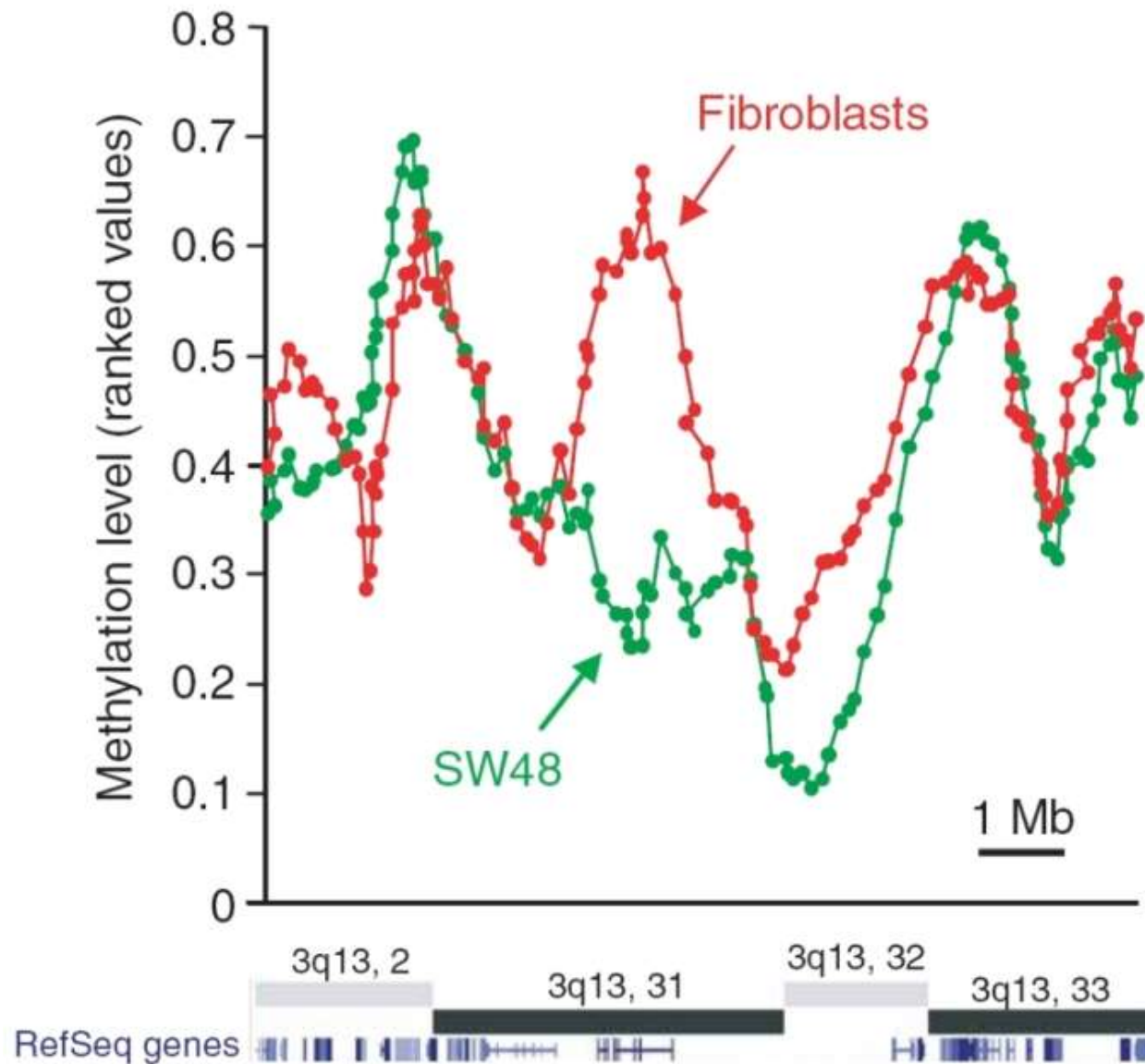
# metylacja DNA



# profilowanie metylacji DNA

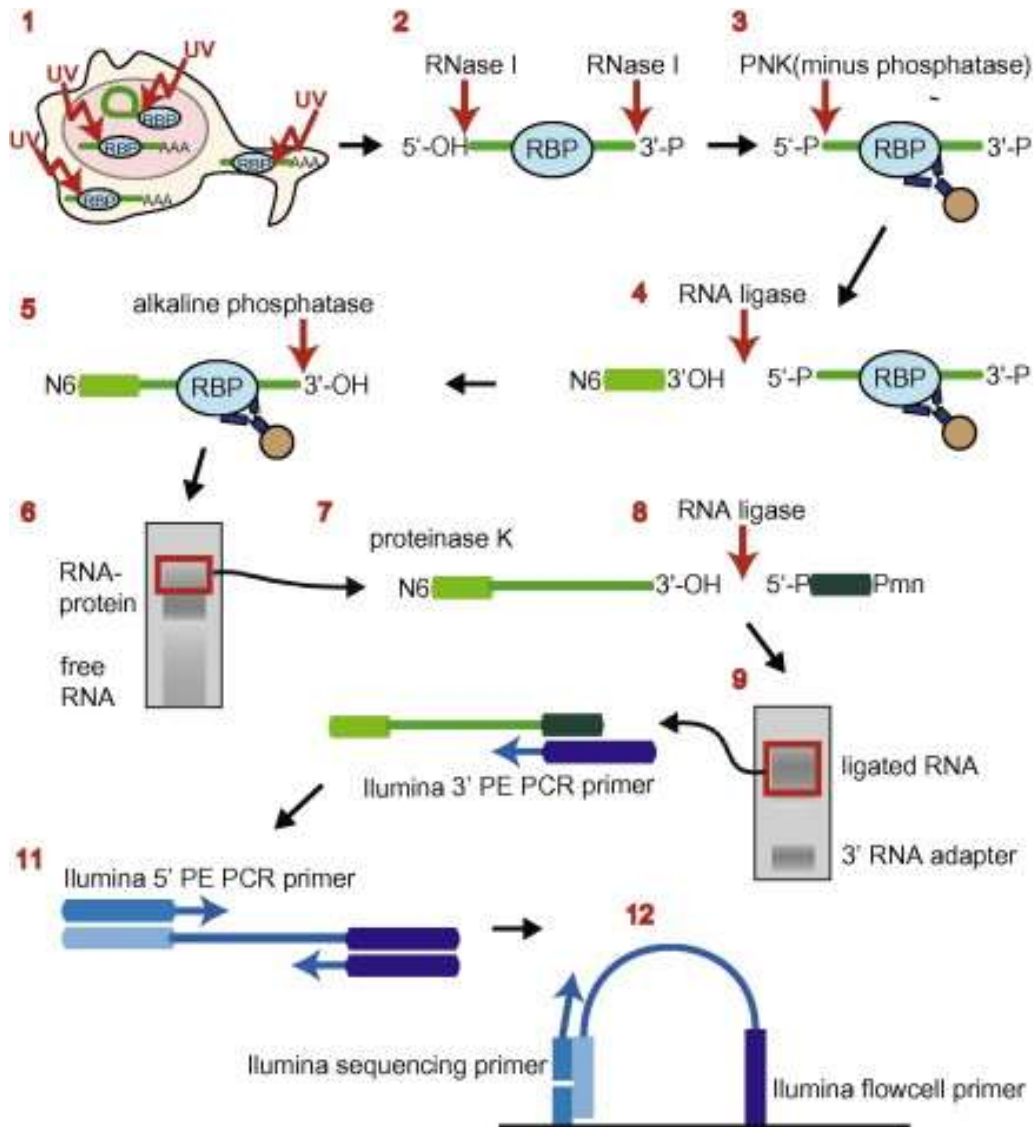


# profilowanie metylacji DNA



# Oddziaływania białek z RNA (CLIP)

## Cross-linking and immunoprecipitation



1. UV crosslink Cells or Tissue.

2. Partial RNA digestion.

3. Immunoprecipitate RBP and phosphorylate RNA 5' end.

4. Ligate the 5' RNA adapter.

5. Dephosphorylate RNA 3' end.

6. Purify RBP-RNA on SDS-PAGE

7. Digest the RBP.

8. Ligate the 3' RNA adapter.

9. Purify RNA on urea-TBE gel.

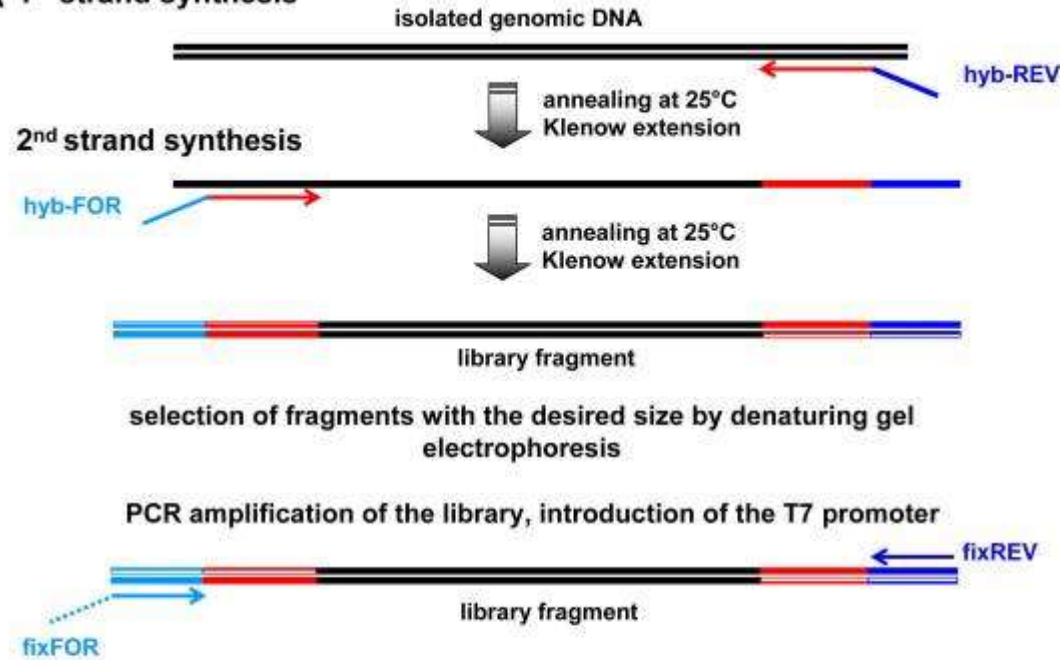
10. Reverse transcription.

11. PCR.

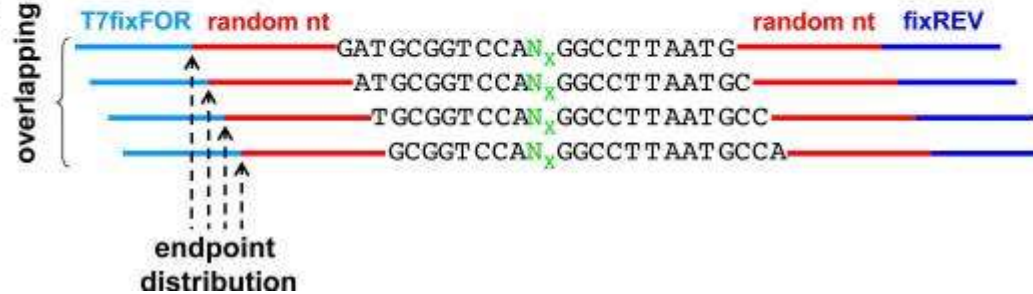
12. Illumina paired-end sequencing.

# Genomic selex

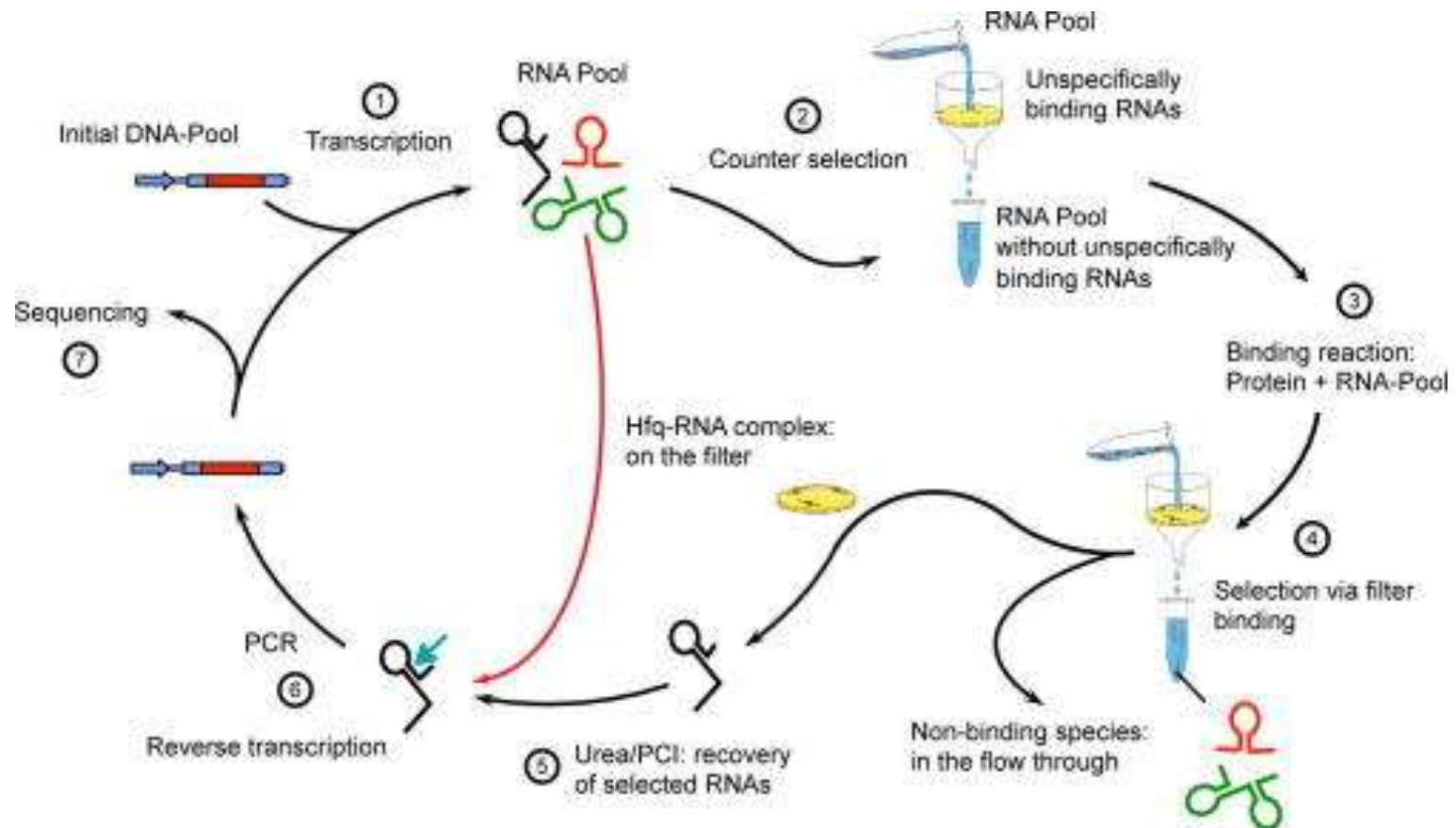
## A 1<sup>st</sup> strand synthesis



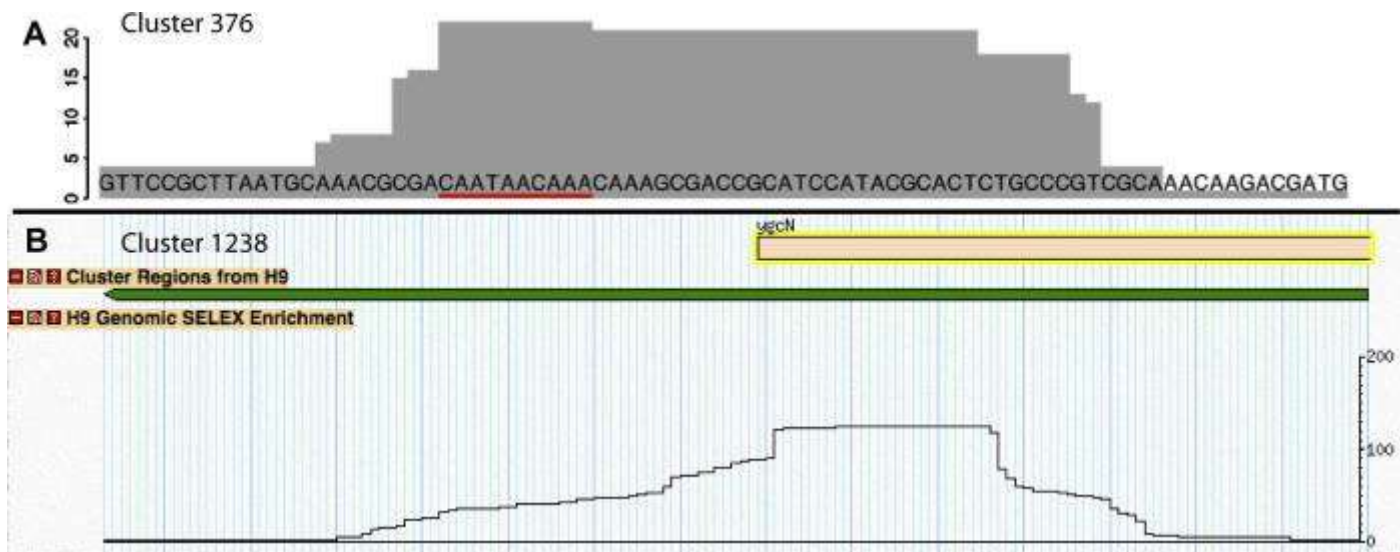
## B



# Genomic selex

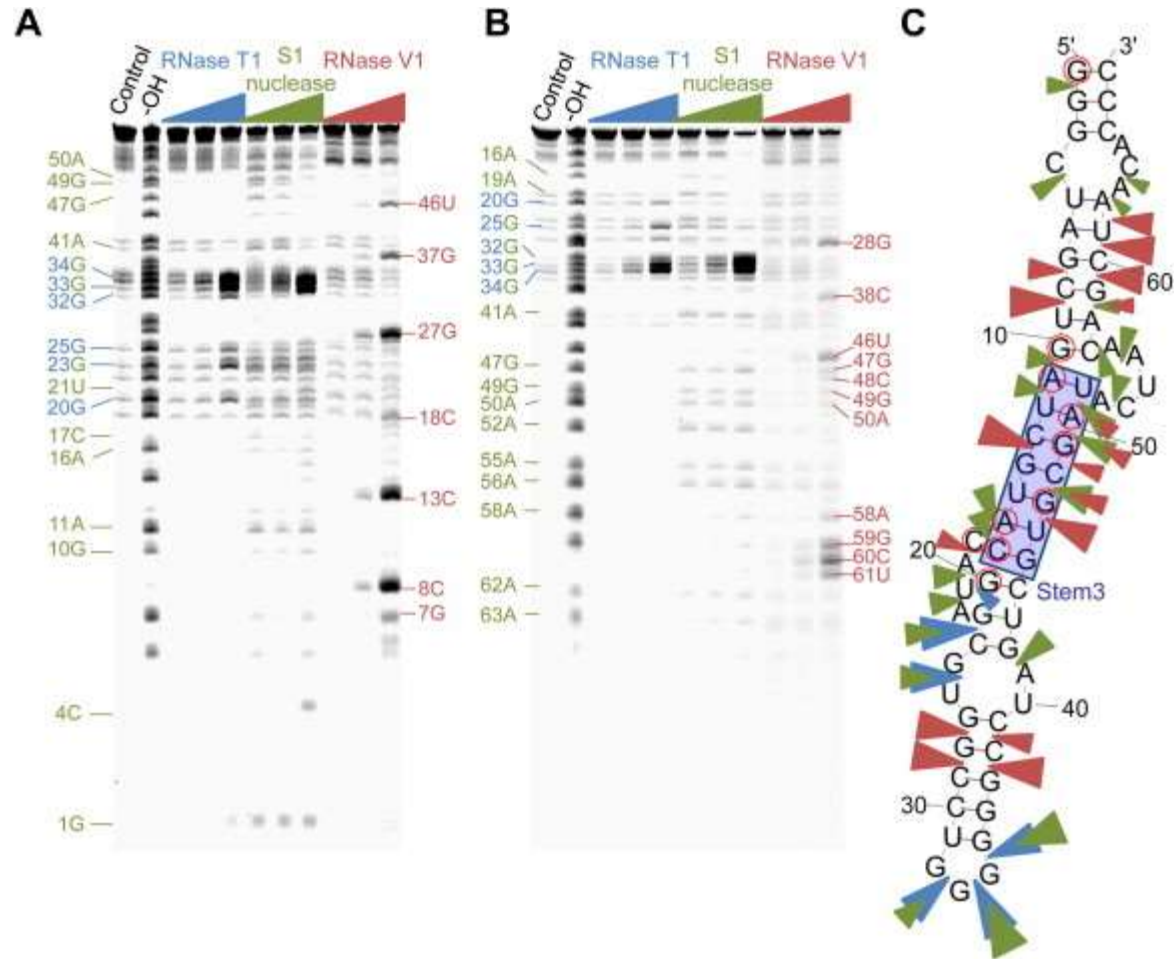


# Genomic selex

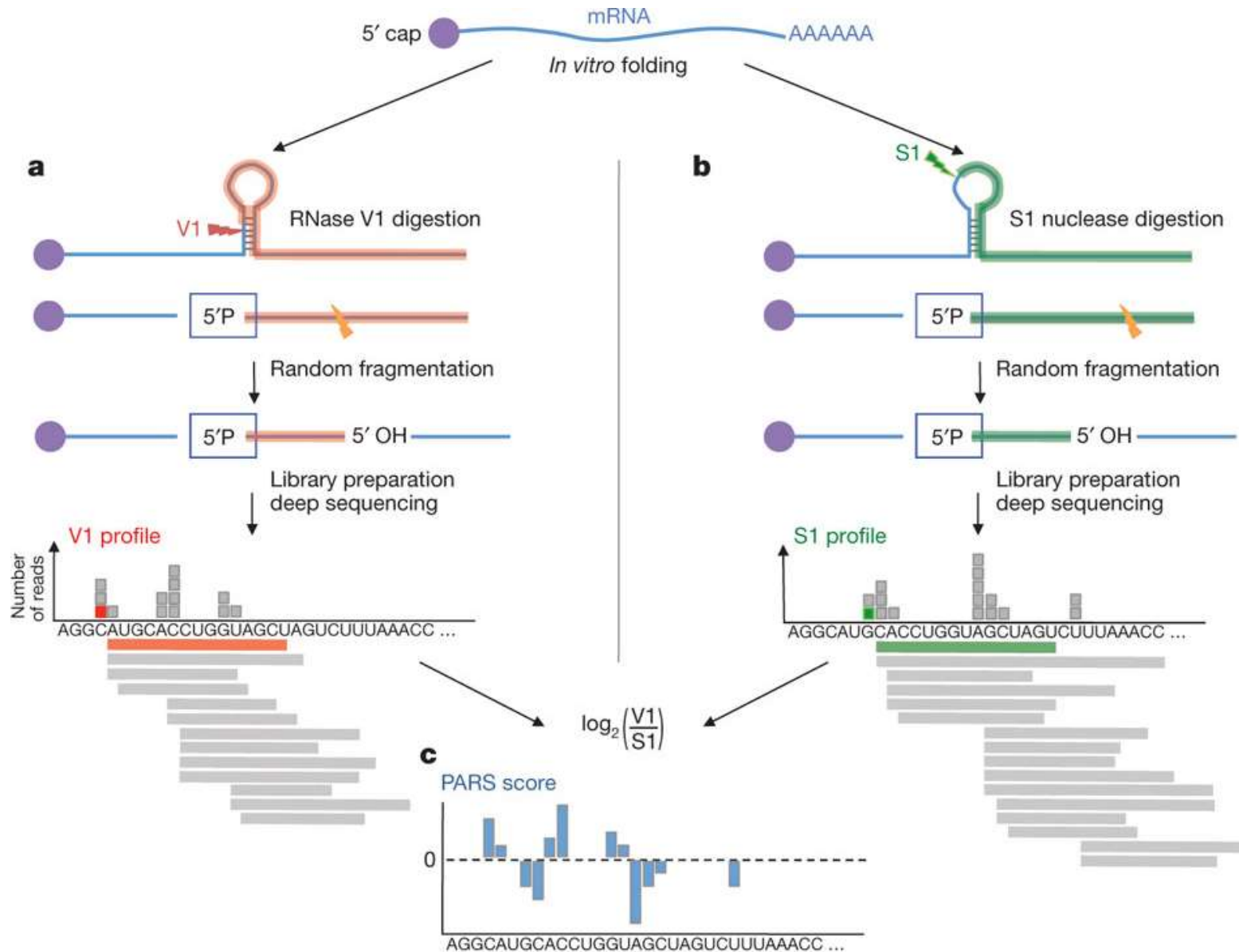




# badanie struktury RNA



# badanie struktury RNA



# badanie struktury RNA

