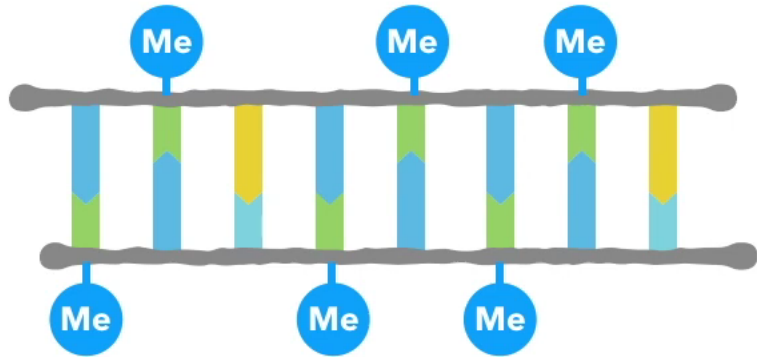# Methylation analysis

## Architecture of large projects in bioinformatics
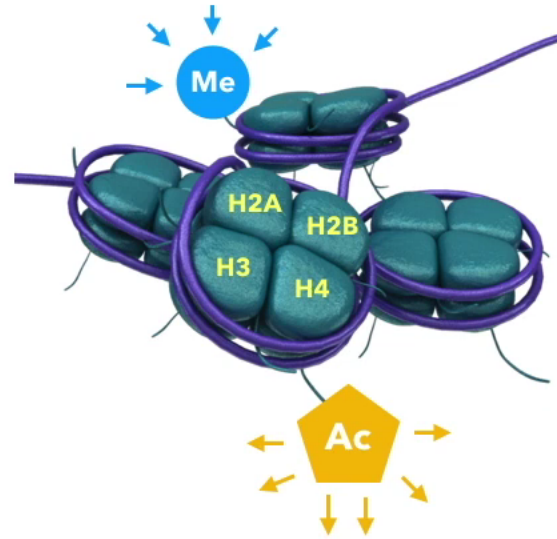
Krzysztof Łukasz

# What is epigenetics

**DNA Methylation**

**Histone Modification**
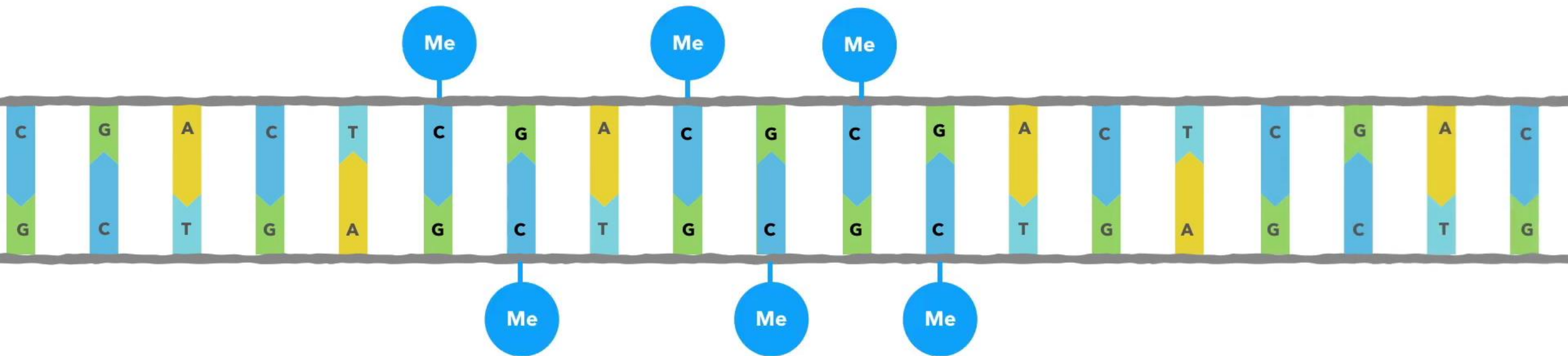
**Non-coding RNA**

Me Me Me
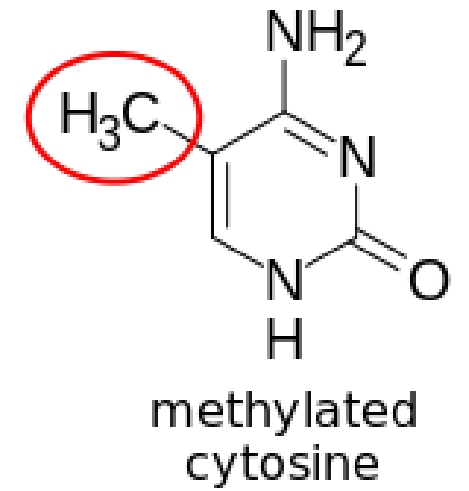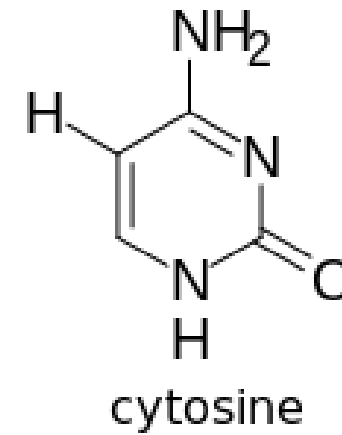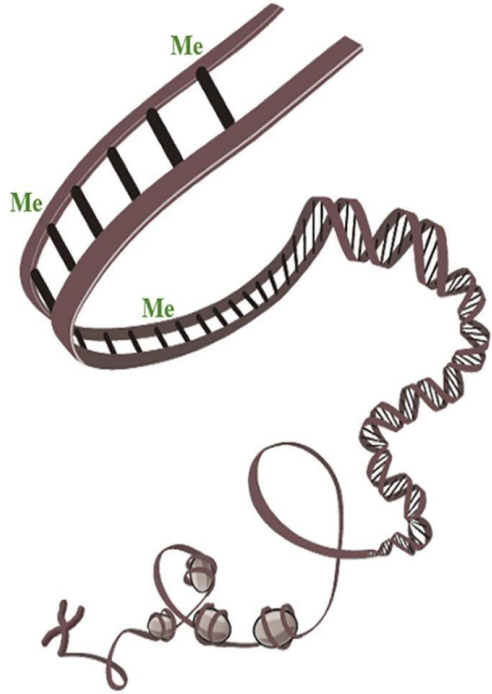
Me Me Me

Me

H2A H2B H3 H4

Ac

*"the branch of biology which studies the causal interactions between genes and their products which bring the phenotype into being."*
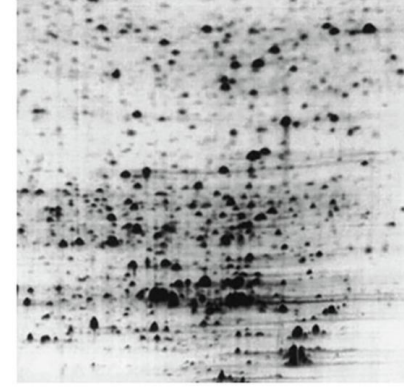
- Conrad Waddington 1940s

# DNA methylation

DNA methylome

Liquid chromatography

Electrophoresis

Microarray

Third-generation sequencing

Second-generation sequencing

First-generation sequencing

Single stranded DNA

Nanopore

A  C  G  T  4mC

Sanger Sequencing

A
T
A
C
G
G
C
T
G
A
C
G
C

Laser  Detector

Capillary electrophoresis

# Epigenetic modification prediction tools based on the type of epigenetic mechanism (information until March 2020).



Bioinformatic tools for DNA methylation and histone modification: A survey, Chenarani et al., Genomics, 2011

# BS experiment design



Input DNA

Bisulfite Conversion

Adaptor Ligation

Amplification

Sequencing

CACATGGTGAAACCCAT
ACATGGTGAACCCATTA

Next Generation Sequencing

# Restrited BS

- Whole-genome methylation – expensive
- Restrictive enzymes
- Analyzing only CpG islands

# Bisulfite-seq workflow

Software survey

# Library preparation

## Directional Library Preparation

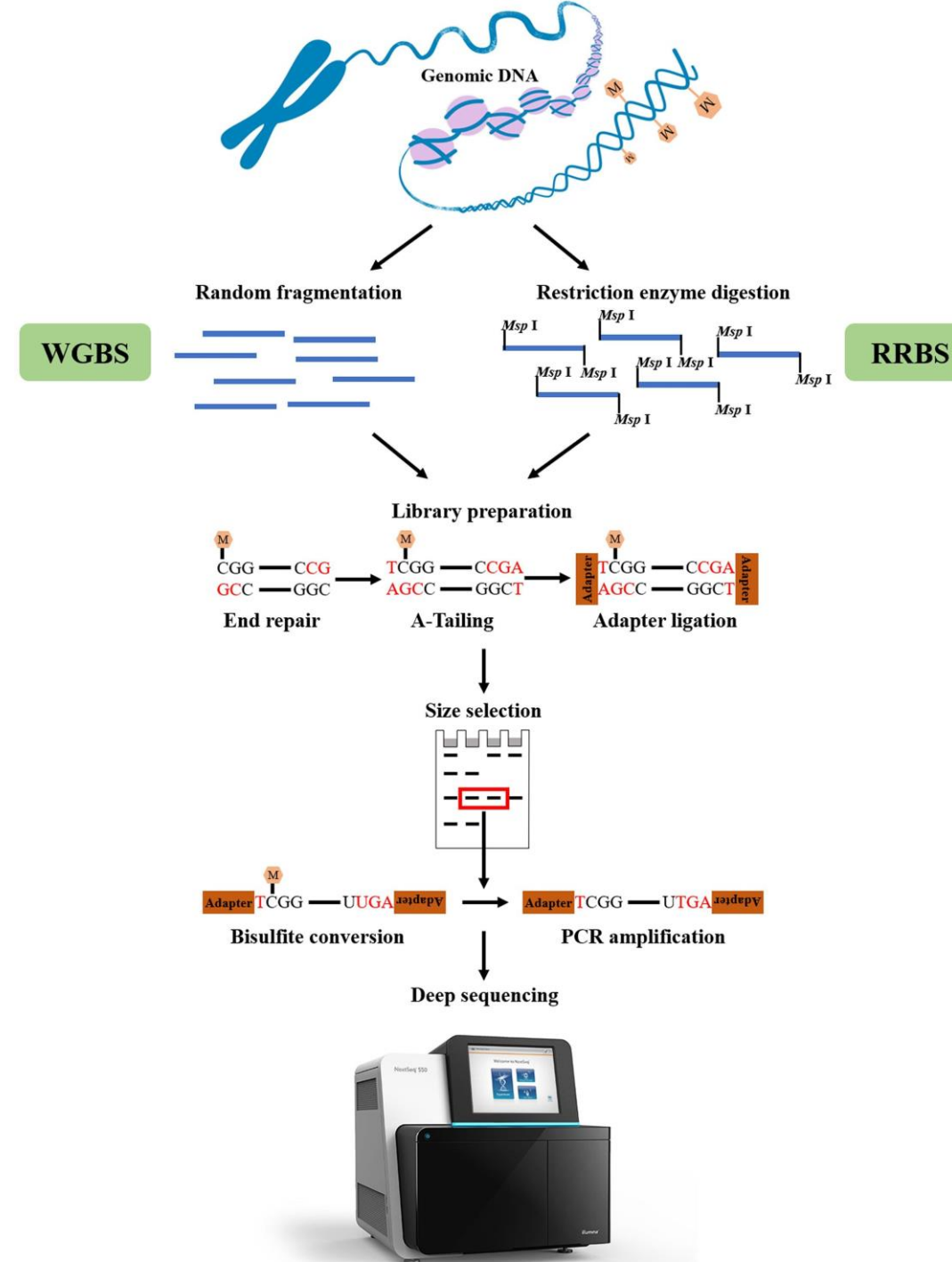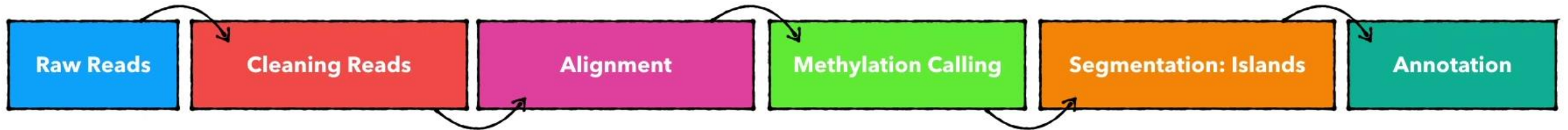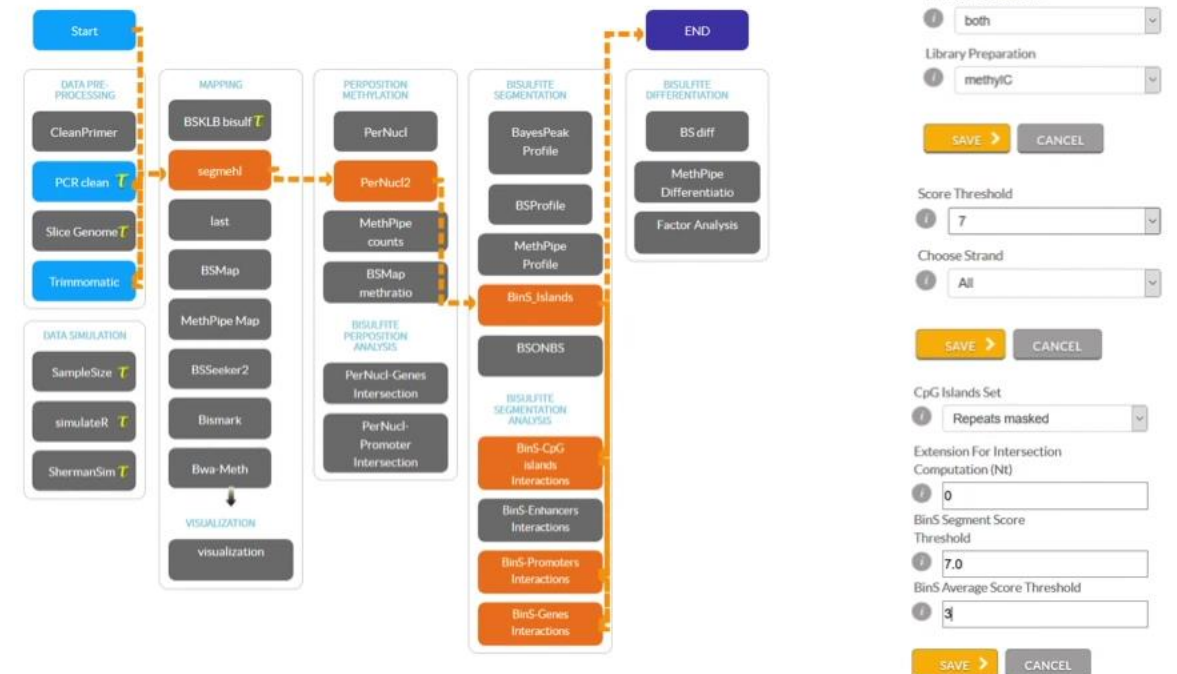5′ …**CC**GG**C**ATGTTTAAA**CGC**T …3′
3′ …GG**CC**GTA**C**AAATTTG**C**GA …5′

bisulfite conversion
PCR amplification

**TT**GG**T**ATGTTTAAA**T**GT**T**          GG**TT**GTA**T**AAATTTG**T**GA

forward strand C → T conversion          reverse strand C → T conversion

## Non-directional Library Preparation

5′ …**CC**GG**C**ATGTTTAAA**C**G**C**T …3′
3′ …GG**CC**GTA**C**AAATTTG**C**GA …5′

bisulfite conversion
PCR amplification

OT      **TT**GG**T**ATGTTTAAA**T**GT**T**
CTOT    **AA**GG**A**ATGTTTAAA**A**GAA**T**

CTOB    GG**AA**GTA**A**AAATTTGA**GA**
OB      GG**TT**GTA**T**AAATTTG**T**GA

forward strand C → T conversion          reverse strand C → T conversion

Non-directional:
- Significantly cheaper
- By far more popular

# Preparing reads

Standard approaches:

• FastQC

• Trimmomatic

• Cutadapt

# BS: Alignment

- Bisulfite treatment introduces mutations into genomic DNA in a methylation dependent manner
- Alignment of BS-seq reads is more challenging
  - Standard alignment methods cannot be used directly
- Standard tools: **Bismarck, BSMap, BWA-Metha**
- Bismark tool uses the following approach to map BS-seq reads
  - Reads from a BS-seq experiment are converted into a C-to-T version **and** a G-to-A version
  - The same conversion for the genome
  - External mapper (Bowtie) alignment to the genome
  - A unique best alignment is determined from four parallel alignment processes

# Bismarck

C-to-T **and** a G-to-A

- **Type**: Command-line tool.
- **Features**: Bismark is widely used for bisulfite sequencing data analysis. It aligns bisulfite-treated sequencing reads to a reference genome and can perform DNA methylation calling.



Figure from (Krueger & Andrews, 2011)

# Methylation calling

- Bismarck, BSMAP and other mappers also perform calling

- Separate software, such as Bis-SNP also available

- For each initial fastq file, we get a call that for each cytosine includes coverage, and methylation frequency

- Now, we need some way to do statistical analysis of such results

A typical methylation call file looks like this:

```
##          chrBase   chr     base strand coverage freqC   freqT
## 1 chr21.9764539 chr21 9764539      R       12 25.00   75.00
## 2 chr21.9764513 chr21 9764513      R       12  0.00  100.00
## 3 chr21.9820622 chr21 9820622      F       13  0.00  100.00
## 4 chr21.9837545 chr21 9837545      F       11  0.00  100.00
## 5 chr21.9849022 chr21 9849022      F      124 72.58   27.42
```

# Segmentation: identification of DMR



**DMR**
Differentially Methylated Region

- Different statistical models
- Fisher's exact test, BSmooth, methylKit, methylSig, DSS, metilene, RADMeth, and Biseq

# Different tools for DMR calling



(a) (b) (c) (d)

Comprehensive Evaluation of Differential Methylation Analysis Methods for Bisulfite Sequencing Data, Int. J. Environ. Res. Public Health 2021, 18(15), 7975;

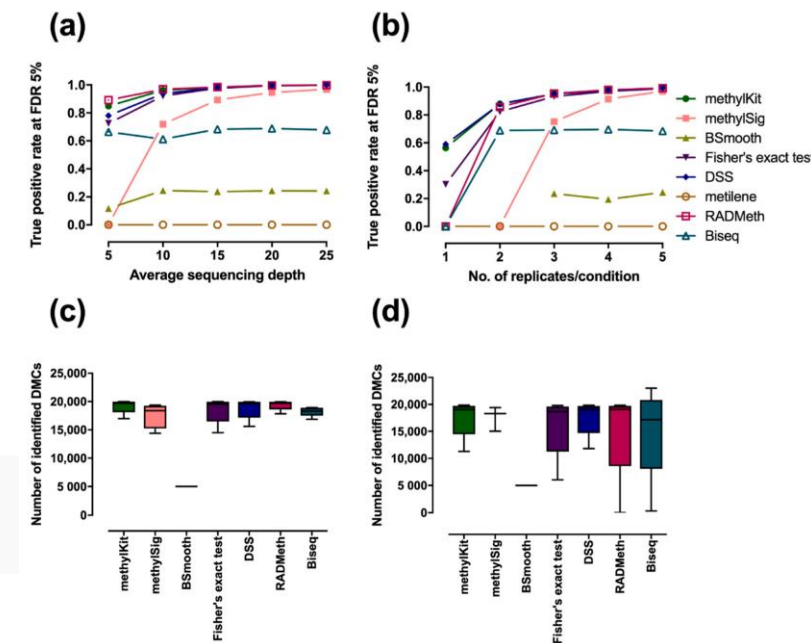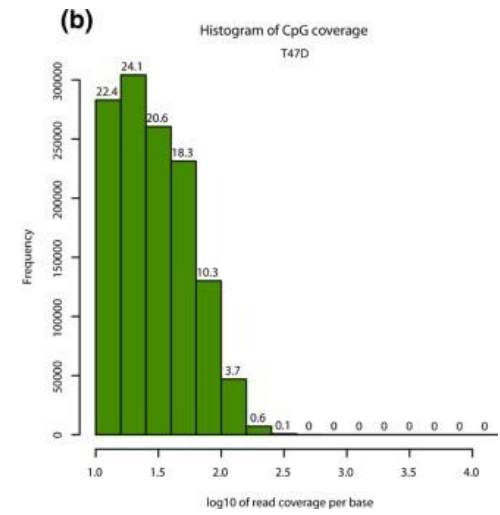| Tool | Version | Model Assumption | Differential Methylation Test | Segmentation | Language | Smoothing |
|------|---------|------------------|-------------------------------|--------------|----------|-----------|
| Fisher's | 1.8.2 | - | Fisher's exact test | tilling window | R | No |
| BSmooth | 1.8.2 | binomial distribution | modified t-test | merging consecutive CpGs | R | Yes |
| methylKit | 0.99.2 | logistic regression | logistic regression test | tilling window or predefined regions | R | No |
| methylSig | 0.4.4 | beta-binomial model | likelihood ratio test | tilling window | R | No |
| DSS | 2.12.0 | Bayesian hierarchical model | Wald test | merging CpGs based on *p*-value | R | No |
| metilene | 0.2–6 | Nonparametric method | 2D Kolmogorov–Smirnov | circular binary segmentation | C | No |
| RADMeth | - | beta-binomial regression | log-likelihood ratio test | correlation between *p*-value pairs within a bin | C++ | No |
| Biseq | 1.12.0 | Beta regression model | Wald test | merging consecutive CpGs | R | Yes |

- Notable variations among methods, and no single method consistently performed best in all benchmarking
- For DMR analysis, methylKit and Fisher's exact test covered more DMRs than other methods

# Methylkit for BS-seq data analysis

# Annotation



- Function of differently methylated genes, comparing with known databases
- "clusterProfiler" is an example of R package to perform:
  - single-gene GO
  - KEGG enrichment.
  - GSEA enrichment analysis

# Annotation



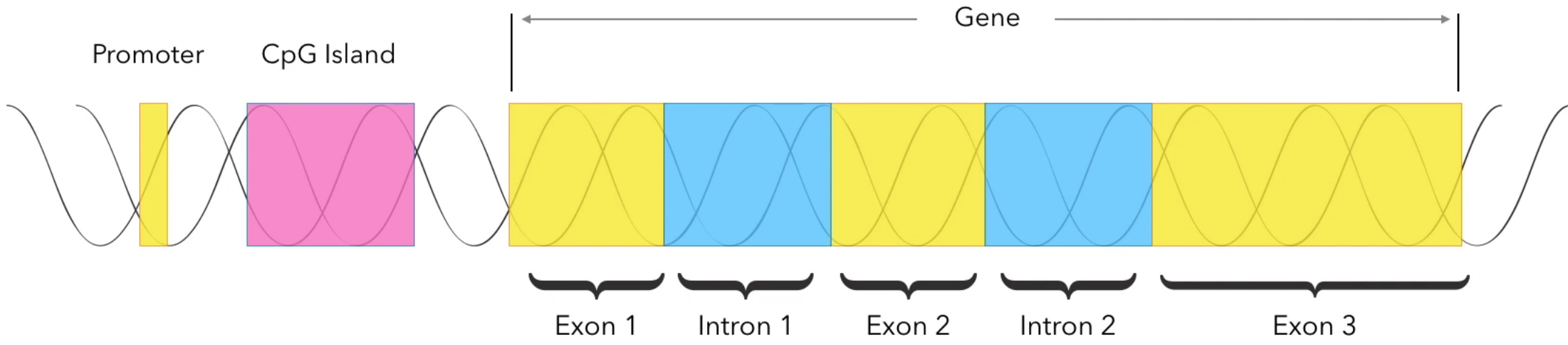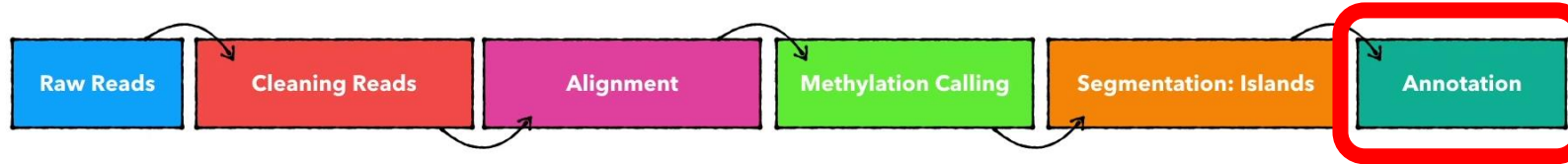| GeneName | GeneID | GeneSrc | Chr | Promoter | Promoter | Promoter | DELIMITE | SegmentStart | SegmentEnd | SegmentAverage | SegmentS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| RP5-857K21.4 | ENSG0000 | HAVANA | chr1 | - | 601436 | 724707 | III | 630909 | 630968 | 3.43038 | 26.5716 |
| AC114498.1 | ENSG0000 | ENSEMBL | chr1 | + | 630896 | 630958 | III | 630909 | 630968 | 3.43038 | 26.5716 |
| RP5-857K21.4 | ENSG0000 | HAVANA | chr1 | - | 601436 | 724707 | III | 634659 | 634674 | 5.67085 | 22.6834 |
| RP5-857K21.11 | ENSG0000 | HAVANA | chr1 | + | 634376 | 634922 | III | 634659 | 634674 | 5.67085 | 22.6834 |
| SKI | ENSG0000 | HAVANA | chr1 | + | 2228695 | 2310119 | III | 2287400 | 2287400 | 10.5384 | 10.5384 |
| NPHP4 | ENSG0000 | HAVANA | chr1 | - | 5862811 | 5992473 | III | 5884629 | 5884629 | 15.5058 | 15.5058 |
| KCNAB2 | ENSG0000 | HAVANA | chr1 | + | 5991466 | 6101193 | III | 6071581 | 6071581 | 15.5063 | 15.5063 |
| KCNAB2 | ENSG0000 | HAVANA | chr1 | + | 5991466 | 6101193 | III | 6071903 | 6071905 | 7.11714 | 12.3272 |
| KCNAB2 | ENSG0000 | HAVANA | chr1 | + | 5991466 | 6101193 | III | 6096719 | 6096720 | 8.80687 | 12.4548 |
| DNAJC11 | ENSG0000 | HAVANA | chr1 | - | 6634168 | 6701924 | III | 6657789 | 6657791 | 8.55709 | 14.8213 |
| CAMTA1 | ENSG0000 | HAVANA | chr1 | + | 6785324 | 7769706 | III | 7222820 | 7222820 | 16.3326 | 16.3326 |
| CAMTA1 | ENSG0000 | HAVANA | chr1 | + | 6785324 | 7769706 | III | 7417151 | 7417154 | 5.9302 | 11.8604 |
| DFFA | ENSG0000 | HAVANA | chr1 | - | 10456522 | 10472526 | III | 10470629 | 10470629 | 11.8334 | 11.8334 |

- Akalin, A., Kormaksson, M., Li, S. *et al.* methylKit: a comprehensive R package for the analysis of genome-wide DNA methylation profiles. *Genome Biol* **13**, R87 (2012). https://doi.org/10.1186/gb-2012-13-10-r87

- Nasibeh C. *et al.,* Bioinformatic tools for DNA methylation and histone modification: A survey, *Genomics*, 2021. https://doi.org/10.1016/j.ygeno.2021.03.004

- Piao, Y.; Xu, W.; Park, K.H.; Ryu, K.H.; Xiang, R. Comprehensive Evaluation of Differential Methylation Analysis Methods for Bisulfite Sequencing Data. *Int. J. Environ. Res. Public Health* **2021**, *18*, 7975. https://doi.org/10.3390/ijerph18157975

- Shizhao Li, Trygve O. Tollefsbol, DNA methylation methods: Global DNA methylation and methylomic analyses, Methods, 2021, https://doi.org/10.1016/j.ymeth.2020.10.002.