OXFORD
UNIVERSITY PRESS

Briefings in Bioinformatics

## DNA methylation analysis in plants: review of computational tools and future perspectives

SCHOLARONE™
Manuscripts

1

2

# DNA methylation analysis in plants: review of computational tools and future perspectives

3    **Jimmy Omony[1,*$], Thomas Nussbaumer[2,3$] and Ruben Gutzat[4*]**

4    [1]Plant Genome and Systems Biology, Helmholtz Center Munich-German Research Center for
5    Environmental Health, 85764 Neuherberg, Germany.

6    [2]Institute of Network Biology, Department of Environmental Science, Helmholtz Center Munich,
7    85764 Neuherberg, Germany.

8    [3]Chair and Institute of Environmental Medicine, UNIKA-T, Technical University of Munich and
9    Helmholtz Center Munich, Research Center for Environmental Health, Augsburg, Germany; CK
10   CARE Christine Kühne Center for Allergy Research and Education, Davos, Switzerland.

11   [4]Gregor Mendel Institute of Molecular Plant Biology, Austrian Academy of Sciences, Vienna
12   BioCenter (VBC), 1030 Vienna, Austria.

13

14   [$] joint first authors


15   **Correspondence*:**

16   Dr. Jimmy Omony, Helmholtz Center Munich, Germany.
17   Email: jimmy.omony@helmholtz-muenchen.de
18
19   Dr. Ruben Gutzat, Gregor Mendel Institute of Molecular Plant Biology, Austria.
20   Email: ruben.gutzat@gmi.oeaw.ac.at
21
22
23

24   **Jimmy Omony** is a postdoc (Bioinformatician) at the Plant Genome and Systems Biology,

25   Helmholtz Center Munich, Germany. His research interests include plant genomics, epigenetics,

26   machine learning, and biostatistics. He undertook the first postdoc at the University of Groningen

27   (RuG). He holds a PhD in computational systems biology (Wageningen University).

28

29   **Thomas Nussbaumer** is a postdoc (Bioinformatician) at the Institute of Network Biology and

30   also in the Institute of Environmental Medicine. His research interests include epigenomics, plant

31   genomics, protein-protein interaction analysis, and microbiomics. He undertook his first postdoc

32   at the University of Vienna and is currently a Postdoc Fellowship Program holder at the

33   Helmholtz Center Munich.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

34

35    **Ruben Gutzat** is a postdoc at the Gregor Mendel Institute of Molecular Plant Biology, Austrian

36    Academy of Sciences, Vienna (Austria). His research interests are in plant epigenetics and

37    developmental biology.

38
39
40
41

## Abstract

Genome-wide DNA methylation studies have quickly expanded due to advances in next-generation sequencing techniques along with a wealth of computational tools to analyze the data. Most of our knowledge about DNA methylation profiles, epigenetic heritability, and the function of DNA methylation in plants derives from the model species *Arabidopsis thaliana.* There are increasingly many studies on DNA methylation in plants – uncovering methylation profiles and explaining variations in different plant tissues. Additionally, DNA methylation comparisons of different plant tissue types and dynamics during development processes are only slowly emerging but are crucial for understanding developmental and regulatory decisions. Translating this knowledge from plant model species to commercial crops could allow the establishment of new varieties with increased stress resilience and improved yield. In this review, we provide an overview of the most commonly applied bioinformatics tools for the analysis of DNA methylation data (particularly bisulfite sequencing data). The performances of a selection of the tools are analyzed for computational time and agreement in predicted methylated sites for *A. thaliana,* which has a smaller genome compared to the hexaploid bread wheat. The performance of the tools was benchmarked on five plant genomes. We give examples of applications of DNA methylation data analysis in crops (with a focus on cereals) and an outlook for future developments for DNA methylation status manipulations and data integration.

**Keywords:** epigenomics, epigenetics, bisulfite sequencing, DNA methylation, plants, differentially methylated regions.

## Introduction

Methylation of cytosine at carbon position 5 (also termed 5-meC) is a hallmark of an epigenetic modification and 5-meC has been described as the 5th base of DNA [1]. Although the extent and context of 5-meC vary considerably between different plant lineages, all plants whose genomes have been sequenced and analyzed so far show substantial DNA methylation [2, 3]. Two major genomic contexts can be distinguished: (i) methylation on gene bodies and (ii) methylation on repeat sequences and transposons. Gene body methylation typically peaks on exons of moderately transcribed genes and, despite a comprehensive body of publications [3-5], its function remains mysterious [6]. Methylation on repeat sequences and transposons is crucial for suppressing transcription and is necessary for establishing heterochromatic domains. Consequently, mutations that abolish most DNA methylation lead to transposon activation and genomic meltdown after several generations in *Arabidopsis thaliana*. However, in early generations, the mutation can be outcrossed and selfed offspring will be isogenic but with different DNA methylation states [7-9]. Experiments along these lines have established that these differences in DNA methylation can be stably inherited over many generations and influence ecologically relevant phenotypic traits [10-15].

In contrast to animals, which only maintain CG methylation, in most plants 5-meC occurs also in several sequence contexts (CG, CHG, and CHH, where H is any of the bases A, T, or C) and is catalyzed by different methyl-transferases acting on different DNA methylation pathways. In *A. thaliana,* CG methylation is maintained by MET1, CHG methylation by CMT3, and CHH by CMT2 and the RNA induced DNA methylation pathway (RdDM). CG methylation occurs in euchromatin and heterochromatin whereas CHG and CHH methylation decorate repeats and transposons [16]. The cross-functioning and redundant DNA methylation pathways form a nuclear/DNA protection system that aids in identifying invading transposons and permanently shutting off their expression (see review by Kim *et al*. [17]).

89   Lister and Ecker [18] argued that 5-meC should be used as a dynamic fifth letter of the genomic

90   code because of the important implications of methylation. It has become tractable to analyze

91   genome-wide DNA methylation states in populations or across different plant species because of

92   advances in next-generation sequencing (NGS) technologies. Much effort has been undertaken

93   to determine the landscape of DNA methylation changes especially in *A. thaliana* and other land

94   plants such as rice and tomato, which have had reference genomes available for several years

95   [19, 20]. DNA methylation patterns vary widely among animals; *Drosophila* completely lacks CG

96   methylation while the human genome is highly methylated (~75% of the cytosines). In *A.*

97   *thaliana*, ~24% of the CGs, ~ 6.7% of the CHGs, and ~1.7% of the CHHs are methylated [21,

98   22].

99   Plants have varying levels of repeat content, which might be the result of bursts of single-repeat

100  retro-elements, which can amplify rapidly using a reverse transcription step to make multiple

101  copies, or DNA transposons, which use a copy-and-paste strategy [23, 24] and thus can amplify

102  during DNA replication. While the repeat content is only ~20% in *Arabidopsis*, in cereals such as

103  barley and wheat the repeat content can be up to 90%. Together with the presence of three

104  subgenomes in hexaploid wheat, these repeats requires tightly regulated epigenetic mechanisms

105  [25]. Genes have evolved different mechanisms for tolerating transposable elements (TEs) in

106  their vicinity [26, 27]. Hirsch and Springer [28] provide a review of the interactions between TEs

107  and gene expression in plants. They discuss three mechanisms by which transposons influence

108  gene expression, namely: (i) the prevailing evidence that TE insertions within introns or

109  untranslated regions of genes are often tolerated and have minimal impact on gene expression

110  levels or splicing. Conversely, TE insertions within genes lead to aberrant or novel transcripts; (ii)

111  TEs act as novel alternative promoters – with the potential to result in different expression

112  patterns; and (iii) TE insertions near genes can influence gene regulation. In Arabidopsis two

113  genes (IBM1 and IBM2) have been identified that prevent spreading of CHG and CHH

114  methylation from transposons into gene bodies or promotors.

115  Interestingly, DNA methylation levels can also affect how plants respond to stress. *Arabidopsis*

116  mutants with reduced global DNA methylation show increased expression of defense related

117  genes and enhanced resistance to pathogens [29]. Polymorphisms of CMT2 correlate with DNA

118  methylation variation along a longitudinal temperature gradient in natural populations [30] and

119  *cmt2* plants are more heat tolerant [31]. Isogenic lines with different DNA methylation states

120  show differences in their ability to compete in synthetic plant communities [32]. Similar influences

121  on stress tolerance have also been observed in monocots, and wheat with experimentally

122  reduced DNA methylation show resilience to salt and oxidative stress. The dynamics of the

123  methylation state of genomic elements are tissue-specific (for instance, in *A. thaliana* seedlings

124  [33-35]) and differ between juvenile and mature plants (e.g. in a study of *Acacia mangium* [36]).

125  Reduced DNA methylation also results in abnormal plant development in *A. thaliana* [37]; hence,

126  an optimally regulated level of methylation is vital for normal plant growth and development.

127  Plant-pathogen invasion can also influence methylation levels in different ways. For instance,

128  genome-wide hypomethylation and hypermethylation influence resistance-related genes [38] and

129  alter gene expression profiles, resulting in plant adaptation to stress. Wang *et al*. [39] showed

130  that drought-induced alterations to DNA methylation in rice influence an epigenetic mechanism

131  that regulates gene expression. As a major modification of the eukaryotic genome, DNA

132  methylation significantly influences gene expression. Methylation of genomic features can lead to

133  different gene regulatory effects. For instance, alteration of a gene's expression potential is a

134  result of DNA methylation affecting the interaction between transcription factors and DNA with

135  chromatin proteins [40]. Additionally, methylation of the promoter region results in repression of

136  gene expression and gene body methylation leads to the opposite effect [41, 42]. Studies have

137  shown that gene body methylated genes are constitutively expressed in a wide range of

138  conditions and tissues [6].

## Chemistry of bisulfite conversion and sequencing

140  Bisulfite sequencing is generally done in three main steps, namely: (i) denaturing, (ii) bisulfite

141  treatment, and (iii) polymerase chain reaction (PCR) amplification. In bisulfite conversion, DNA is

142  denatured in a process that separates the forward and reverse strands. This is followed by

143  treatment with sodium bisulfite, which converts unmethylated cytosine into uracil – which is then

144  converted to thymine during PCR [43]. Quantification of the abundance of each cytosine can be

145  achieved via Sanger sequencing [44] or NGS technologies [45]. The DNA strands cease to be

146  complementary after bisulfite conversion. Treatment of genomic DNA with sodium bisulfite [46]

147  enables us to distinguish between highly similar (and yet different) methylated cytosine, which

148  has the same base-pairing features as unmethylated cytosine. Mapping read sequences to a

149  reference genome enables the determination of positions with matching and mismatching bases.

150  This process enables identification of methylated and unmethylated bases.

151  Bisulfite sequencing can be accomplished with different sequencing kits depending on whether

152  whole-genome bisulfite sequencing (WGBS) or reduced representation bisulfite sequencing

153  (RRBS) (WGBS: Lister and Ecker [18], RRBS: Jeddeloh *et al*. [47], Schmidt *et al*. [48]) is

154  performed. Currently, WGBS remains the most informative method for generating DNA

155  methylation data. It provides a huge wealth of data and requires no prior targeting. Unlike

156  WGBS, which is expensive, RRBS can be performed more economically because it is restricted

157  to CpG-enriched regions that make up a smaller portion of the genome. The restriction enzyme

158  *Msp1* cleaves at 5'-C*CGG-3' targets (base preceding * is methylated), thereby, mainly CpG-rich

159  regions are targeted – which is advantageous for large genomes.

## Typical workflow for processing bisulfite sequencing data

161  Before reads are mapped to a reference genome, the sequencing quality of reads can be

162  checked with FastQC (http://www.bioinformatics.babraham.ac.uk/projects/fastqc) or NGS QC

163  Toolkit [49] followed by removing low-quality bases and adapters with, among others, Trim

164  Galore    (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/),    cutadapt    [50],    or

165  Trimmomatic [51]. However, some WGBS data processing tools integrate various analytic steps -

166  enabling data preprocessing, read alignment, a more robust statistical analysis which output

167  statistics such as read coverage, the percentage of uniquely aligned reads, and statistics on the

168  three methylation contexts (CpG/CHG/CHH). One such tool is gemBS [52], which is a recently

169  published pipeline for processing and analysis of WGBS data. The pipeline integrates data pre-

170  processing and analysis steps from adaptor trimming through downstream statistical analysis of

171  mapping results. gemBS uses the high-performance read aligner GEM3 [53] as a dependency

172  and BScall (embedded in samtools, bcftools; http://samtools.sourceforge.net/) which is a variant

173  caller for bisulfite sequencing data. Both GEM3 and BScall support single and paired-end reads.

174  Further reading on the generic workflow of analyzing WGBS is found in the work of Liang *et al*.

175  [54] and Wrecyzcka *et al*. [55].

## Non-bisulfite based methods and related bioinformatics tools

177  While bisulfite sequencing methods represent the most popular approaches for analyzing

178  epigenomic data, there are other approaches within the field of DNA modification based

179  methods. These approaches include MeDIP-seq and MethylCap-seq, in methylated DNA

180  immunoprecipitation (MeDIP) analyses [56] where the genomic DNA is randomly sheared,

181  sonicated, and immunoprecipitated with an antibody recognizing 5-methylcytidine. Precipitated

182  DNA can either be sequenced or hybridized to microarrays. MethylCap-seq uses the methyl-CpG

183  Binding Domain (MBD) of MeCP2 [57] while oxBS [58] is used to specifically detect 5-

184  methylcysteine (5mC) and 5-hydroxymethylcytosine (5hmC) which can be also done with *Tet*-

185  assisted bisulfite sequencing (Tet) [59]. CAB and fCAB for the recognition of 5caC [60]. Notably,

186    the presence/absence of 5hmC in plants remains contentious. Some scholars claim that 5hmC is

187    present in plants [61, 62] while others claim its absent [63]. A comprehensive overview of the

188    various tools is given at https://omictools.com/medip-seq-category.

## Tools for analyzing epigenomics datasets

190    Bismark [64] and BSMap [65], as one of the first published tools for quantifying epigenomic

191    datasets had to address the challenge of attaining high read mapping efficiency to enable a

192    sensitive sequence search. Bowtie [66], Merman [67], SNAP (http://snap.cs.berkeley.edu/), and

193    Bowtie2 [68] have been used as dependencies in epigenomics tools, for instance, BS-Seeker

194    [69], BS-Seeker2 [70], BS-Seeker3 [71], BRAT-nova [72], WALT [73] and Bismark, which are

195    currently among the most commonly applied tools for mapping bisulfite methylation data. We

196    outlined the most common tools for mapping bisulfite sequencing data along with tools that allow

197    for the detection and analysis of differentially methylated regions (DMRs). The program

198    parameters as well as input and output data formats are specified in Table S1. This table

199    provides an overview of the main tools for mapping and analysis of epigenomic data –

200    particularly for bisulfite sequencing data. Additionally, we also categorized the tools into three

201    major classes, namely: (a) mapping, (b) statistical analysis, and (c) complete pipelines (Table

202    S1). The defining features for each tool, such as their ability to handle single or double-stranded

203    sequence data as well as their ability to process data and perform down-stream statistical

204    analysis, are also provided. Reviews by Adusulalli *et al*. [74], Shafi *et al*. [75] and Wrecyzcka *et*

205    *al*. [55] complement our overview Table S1. The most frequently applied computational

206    epigenetics methods were applied and tested using DNA methylation data, particularly with data

207    acquired from bisulfite sequencing experiments. Therefore, there are many statistical procedures

208    available for analyzing methylome data – categorized into the parametric and non-parametric

209    approach. Both approaches are widely used in the literature [76]. For instance, MethylMix [77] is

210    an excellent example of a parametric approach which uses Bayesian mixture models to identify

211  DNA methylation states of genes as either hypo- or hypermethylated. The method entails fitting a

212  distribution function onto the frequencies of DNA methylation counts. The advantage of using

213  non-parametric models is that no prior knowledge of the data distribution is required. However,

214  when such knowledge is available, then parametric models are the preferred choice for

215  modelling such data. MethylMix quantifies the effect of DNA methylation on genes, which is

216  interesting for integrative studies that aim at establishing the association between the

217  methylation states of the individual genes and their expression profiles. Investigating such

218  associations unravels any hidden variations within and between samples (or tissues) as

219  illustrated in [78-80]. Lea *et al*. [81] discussed the applications of mixed models on DNA

220  methylation in plant epigenetics. They specifically focused on the binomial mixed model with the

221  sampling-based algorithm (MACAU: Mixed model association for count data via data

222  augmentation) for the approximation of parameters and computation of *p*-values. Other

223  modelling frameworks are based on algorithms that integrate various analytical steps resulting in

224  the detection of DMRs across the entire genome, for instance: (i) the weighted optimization

225  algorithm proposed in [82] (which is an extension of MethylKit [83]), and (ii) ChAMP.DMR [84]

226  which applies the Bumphunter [85] or ProbeLasso Algorithm [86]. An example of a non-

227  parametric model is the Bayesian approach based on the Dirichlet-process beta-mixture model –

228  which is used for clustering methylation profiles [76]. The model considers the DNA methylation

229  expressions consisting of an infinite number of beta mixture distributions [87, 88].


## DNA methylation: plant physiology and pathophysiology

231  Investigating the dynamics of DNA methylation in plant growth and development requires the

232  analysis of samples from different plant tissues (e.g. Bartels *et al*. [34]). To our knowledge, no

233  existing software has been developed specifically for the analysis of plant physiology and

234  pathophysiology. However, there are many studies analyzing bisulfite data using samples from

235  different plant developmental stages (from seedlings to mature plants). For instance, Bismark –

236  in leaf tissues from bread wheat seedlings [89], BSMap – for various datasets from different

237  tissues in *A. thaliana* [90], and BS-Seeker2 – for young *Zea mays* leaves [91]. With rapid

238  advancements in the development of software/tools for analysis of epigenomes, we are

239  optimistic such tools will soon be available to the public.

## Differentially methylated regions and their significance

241  Genomic regions (or bases) with different methylation profiles between samples are known as

242  differentially methylated regions (DMRs). This is also referred to as differentially methylated CpG

243  sites since the CpG-methylated sites occur in much larger numbers compared to the non-CpG

244  contexts (CHG and CHH) [92, 93]. Peak detection enables the identification of CpG islands –

245  which are essential for differentiating methylation profiles between samples (typically between

246  controls and test samples). CpG islands are not randomly distributed in the genome but are

247  instead grouped close together [94]. Long stretches of non-dense CpG sites, known as CpG

248  shores can also be detected. Combining the methylation profiles of both CpG-islands and CpG-

249  shores enables more efficient comparative analysis of DNA methylation profiles between

250  samples.

251  Various statistical algorithms have been proposed for identifying DMRs – the most popular ones

252  being: methylKit [83], metilene [95], DMRcaller [96], and Bumphunter [85]. For elaborate

253  discussions on the DMR detection methods and a discussion on choosing the right method for

254  DMR detection see Hebestreit *et al*. [97], and Kurdyukov and Bullock [98]. The tools are written

255  and compiled in different programming languages (e.g. R, Python, Perl, Java, C, and C++; Table

256  S1). Essentially, such tools are used to identify DMRs from either targeted regions of the

257  genome or from the whole genome. Critical considerations have to be made, e.g. the choice of

258  experimental designs for experiments and statistical methods for data analysis [99]. DMRs are

259  intricately linked to transcription and the abundance of CpG sites (CpG islands). A high

260  concentration of CpG sites are often found within the promoter regions of genes – so it is

261  essential to accurately identify such sites. Methylation of promoter regions influences the level of

262  transcription – heavy methylation disrupts transcription and de-methylation leads to transcription

263  reactivation [100-102].

264  Peak identification and normalization are crucial initial steps in analyzing DNA methylation data

265  and visualization and can be useful for comparing datasets and judging the performance and

266  agreement between tools. Post-processing and visualization of (differentially) methylated sites

267  enable high-resolution exploration and comparison of regions in the genome for variations in

268  methylation profiles. Therefore, tools like BiQ [103] and BSeQC [104] have aided quality control

269  and visualization of methylation data, thereby enabling researchers to explore data attributes and

270  perform data quality control before analysis. There are many methods for clustering methylation

271  marks such as the dynamic genome warping [105] approach which uses hierarchical clustering

272  and the combination of different epigenomics analytic platforms and data integrative modules.

273  Dynamic genome warping has been demonstrated to be a reliable way to get more meaningful

274  results from datasets (for instance, Chari *et al*. [106]). To utilize this method, Liang *et al*. [54]

275  developed a web-server to analyze whole-genome bisulfite sequencing data and their platform

276  includes major steps for detection and mapping of the conversion rate, detection of DMRs, and

277  their association to gene expression. Wreczycka *et al*. [55] discussed data requirements and

278  computational attributes for specific software and assess bisulfite sequencing data analysis

279  methods, alignment and data processing, detection of differential methylation, and assess

280  strategies for handling large epigenetic datasets. In contrast, our work highlights existing

281  asymmetries between mapping tools and contrasts their computational capabilities.

282  Another important aspect in plant epigenetics is how hypomethylation and hypermethylation

283  affects gene expression. The concept of hypomethylation and hypermethylation is not limited to

284  plants as they have also been extensively studied in cancer progression in humans [107],

285   coronary heart disease [108] and eukaryotes in general [109]. The division of DMRs into hypo-

286   and hypermethylated enables investigations into the influence of both types of methylation on

287   gene expression. Many computational tools have integrated modules that enable the extraction

288   and quantification of the extent of hypo- and hypermethylation in genes. One such tool is

289   MethylMix, which requires that changes in a gene's methylation state must also agree with its

290   expression profile. Additionally, it requires a treatment and control sample (for agricultural

291   studies) or healthy and disease conditions (for clinical studies).

## Downstream analyses of bisulfite methylome data

293   After data processing and calling of methylation sites, downstream analysis can be performed –

294   including the functional annotation of differentially methylated regions and analysis of the

295   associated pathways influenced by the targeted genes. Such analysis enables the assignment of

296   functions and gene annotation as seen in the overviews of Bioinformatics omicX tools

297   (https://omictools.com/epigenomics-category). Examples of tools for performing downstream

298   analysis are given in Table 1.

## Technical challenges: conversion rate, repetitive regions and differentially methylated regions (DMRs)

301   The main challenges in the analysis of DNA methylation data include incomplete methylation

302   patterns and overdispersion of read-mappings [110-112]. Here, overdispersion means the

303   presence of variability in the reads compared to the expected read distributions based on a given

304   model structure. When epigenomics marks coincide with repetitive regions in the genome,

305   mapping tools need to keep reads that map to multiple genomic locations – making these tools

306   slower and computationally memory-intensive. This problem can be partly circumvented through

307   parallel computing using multiple threads, especially for larger repetitive plant genomes.

## Conversion rates

As a method for studying DNA methylation, bisulfite conversion involves the conversion of cytosine to uracil (while 5-methylcytosine, 5-mC remains unchanged). Bisulfite sequence conversion rates vary for different datasets. It is essential for conversion rates to be determined accurately to ensure the reliability of down-stream data analysis. Reliable results can be obtained from datasets with bisulfite-conversion rates higher than ~0.999 (see e.g. Sun *et al*. [113] – demonstrated using their tool MethQA). However, they urge caution for datasets with lower conversion rates. Modern commercially available bisulfite sequence conversion kits generally indicate conversion efficiencies of 90–100% [114]. An elaborate discussion on methods for estimating conversion rate from bisulfite DNA methylation data is provided in [115, 116].

## Description of experiment: benchmarking selected tools

We aimed to determine how the well-established computational epigenomics methods perform on a small genome such as *A. thaliana* with ~130 Mbp (TAIR10) compared to a genome with a high repeat content and much larger genome size such as bread wheat – taking chromosome 1A (Chr1A) for demonstration purpose, IWGSC.v1 *et al*. [117]. We used bisulfite sequencing data from two studies (with accession numbers SRR429549 [118, 119] for *A. thaliana* and ERR1141918 [89] for *T. aestivum*, data from NCBI) and applied four methods: BSMap [65], Bismark [64], BS-Seeker3, and segemehl [120]. Our analysis focused on the speed and agreement of common methylated sites between the tools. BS-Seeker3 was the fastest, followed by BSMap, while Bismark and segemehl were the slowest irrespective of genome size – especially for multiple threads (Figure 1: A and B). When using a single thread, segemehl (keeping reads that mapped a maximum of 3 times) performed slowest compared with the other methods. Overall, the computation time required for the *T. aestivum* (Chr1A) dataset is significantly longer than those from *A. thaliana* (Figure 1: A and B). When comparing the

332    reported sites, we found that, for *A. thaliana*, 562,051 sites are shared amongst all four tools.

333    While most sites were overlapping between BSMap, BS-Seeker3 and Bismark, likely because

334    they use the same mapping software, segemehl reported only ~10% of these sites. However, for

335    *T. aestivum*, ~101,944 sites were reported with most of them being reported in segemehl (Figure

336    1: C and D). The existence of such asymmetries requires more attention and is certainly worth

337    taking into consideration when using the different computational tools. Other studies on

338    comparisons of the performance of epigenetics analysis tools, specifically focusing on mapping

339    short reads for bisulfite sequencing data, can be found in the work of Tran *et al*. [121]. Several

340    studies have also compared run-time and memory consumption of different epigenomics tools,

341    such as Tran *et al.* [121] who compared the five bisulfite short read mapping tools BSMap,

342    Bismark, BS-Seeker, BiSS and BRAT-BW and Bismark performed best on real data, followed by

343    BiSS, BSMap and BRAT-BW and BS-Seeker. Recently, Huang *et al.* [71] proposed BS-Seeker3

344    – a fast mapping tool for bisulfite data, and compared it performance for run-time and sensitivity

345    to sister tools like Bismark, BRAT-nova, and BSMap. Additional to being accurate and versatile,

346    Huang *et al.* concluded that BS-Seeker3 is an ultra-fast pipeline to process bisulfite-converted

347    reads. The tool also aids visualization of methylation data; hence, justifying its comparability to

348    the other three tools (Bismark, BRAT-nova and BSMap).

349    We simulated reads from *A. thaliana* and bread wheat using the tool by Sherman

350    (https://www.bioinformatics.babraham.ac.uk/projects/sherman/) to test the performances of the

351    four tools by comparing the precision and sensitivity along all chromosomes (Figure 2). The

352    sensitivity, also sometimes referred to as recall, is defined as TP/(TP+FN). The precision is

353    defined as TP/(TP+FP), where TP – true positive, FN – false negative and FP – false positive.

354    We observed best performances for the Bismark, followed by BSMap and segemehl, while BS-

355    Seeker3 seemed to have a lower sensitivity in *A. thaliana* compared to the other tools. For bread

356    wheat a similar order to performances of tools was observed when reads where simulated for

357  each subgenomes of chromosome 1 with the three genome copies. All scripts were provided in

358  GitHub (https://github.com/jomony/EPItools/blob/master/README.md).

## Feature comparison between the tools and related literature benchmarking

To further benchmark the performance of the tools, we used bisulfite sequencing data from five plant genomes. These genomes consist of the dicots: *Arabidopsis thaliana* (genome size ~0.13Gb, SRR4295494), *Arabidopsis lyrata* (~0.21Gb, SRR3880297) and *Glycine max* (~1.2Gb, SRR5079790), and also the monocots: *Triticum aestivum* (chromosome 1A, size ~0.67Gb, ERR1141918) and *Oryza sativa* (~0.43Gb, SRR7265433). Figure 3(A) shows the results of a comparative analysis of the memory footprint analysis of the performance of the four tools benchmarked using data from five genomes. These results come from mapping the bisulfite reads data to their respective reference genomes. Association analysis was performed for each of the four tools as seen in the linear regression model fits (Figure 3: B to E). The results show that the genome sizes for each of the five genomes are significantly correlated to the memory footprint analysis ($p$-values < 0.05).

The key attributes and parameters for the four tools are summarized in Table S2. This table presents a summary of the supported features in the four tools (BSMap, BS-Seeker3, Bismark, and segemehl). Such features are essential for deciding on which tool to use for mapping reads and data analysis. Examples of such features can also be found in the work of Guo *et al*. [70] and Tran *et al*. [121]. Lee *et al*. [122] evaluated the mapping accuracy and mapping rates for Bismark, BSMap, and BS-Seeker2 as a function of the error rates. Using whole genome bisulfite sequencing data, they assessed the influence of the error rates on the mapping rates and mapping accuracy and observed that at low error rates (<4%), BSMap had a higher mapping rate than Bismark and BS-Seeker2. On the contrary, BSMap had a lower mapping accuracy than

380  Bismark and BS-Seeker2. They also showed that mapping accuracy is independent of the

381  methylation level.

382  A discussion on benchmarking approaches with a focus on short sequence mapping tools is

383  found in the work of Hatem *et al*. [123]. They assess the performance of various aligners for the

384  read mapping tools and benchmark them using criteria such as mapping percentage, running

385  time and memory footprint. Variations in parameters such as seed length, base quality, single- or

386  paired-end reads on the mapping quality are also evaluated. Benchmarking of tools by

387  comparing the performance of each tool based on multiple attributes can be achieved in various

388  ways, for instance, by assessing the: (i) effect of the read length and sequencing error, (ii) effect

389  of data processing, and (iii) effect of varying parameters in the tools. These are some of the

390  approaches discussed by Tran *et al*. [121]. They compared the performance of epigenomic

391  mapping tools such as BSMap, Bismark, BS-Seeker, BRAT-BW [124] and the Bisulfite

392  Sequencing Scorer (BiSS) [125]. Tran *et al.* primarily benchmarked the performance of the tools

393  basing on mapping efficiency (as the percentage of reads that map uniquely to the genome) and

394  the CPU time.

## Outlook

396  In the near future, there is a need for more comparative analyses to explore the epigenomes of

397  diverse plants in different development stages together with various stress factors. This would

398  enable the discovery of exclusive and common epigenetic regulatory mechanisms. Uncovering

399  and exploiting such mechanisms could potentially promote adaptation to changing environmental

400  conditions. Moreover, a large number of methylomes are required to study the effect of the

401  environment and stress conditions on the epigenomic state of a single plant [126, 127].

402  Resources like the 1001 epigenomes project (https://1001genomes.org/) in *A. thaliana* are

403    exciting datasets to aid in our understanding of the role of the epigenome. However, it remains

404    unclear whether the observations in these studies are directly applicable to crops.

405    Computational tools are instrumental for bridging the gap between mapping of sequenced reads,

406    the accurate prediction of methylated sites, and their statistical analysis However, this effort is

407    hampered by variations in the size of epigenomic marks and the complexity associated with

408    normalizing peaks. The need to increase crop yield on the same amount, and in some cases

409    dwindling, of arable land is another important aspect that requires advancements in epigenomics

410    studies. Several studies have shown that during seed and grain development, the plant

411    epigenome changes and leads to gene silencing. Therefore, a change in the epigenetic state of a

412    plant would result in an increase in its likelihood of adapting from one geographical location to

413    another or to different environmental conditions.

414    Lämke and Bäur [128] argued that such modifications have the potential to provide a mechanistic

415    basis for stress memory in plants. This enables plants to respond more efficiently to recurring

416    stress from the environment, for instance drought and salinity stress [129], a topic that was

417    reviewed by Golldack *et al*. [130] (and more recently by Yang and Guo [131] and Abhinandan *et*

418    *al*. [132]). This might enable plants to prepare their offspring for future attacks from stressors and

419    to improve their adaptation to specific stress factors [130]. Plant adaptation to stress might

420    enable us to explore new ways to improve yield, for instance by shortening or prolonging the time

421    for grain development, by finding ways to regulate the expression of the three homeologs in

422    wheat, or by interfering with fruit ripening (as seen in tomatoes [133-135] and other fruits like

423    peach, apples, and strawberries [136]). A more intriguing discussion on the epigenetic

424    mechanisms of plant stress response and adaptation to different environmental conditions was

425    reviewed in [137-139].

426    In this review, we have discussed the use of bioinformatics tools to study DNA methylation data

427    in plants. Notably, several studies in humans and mouse were successfully performed using

428  popular tools like BSMap, BS-Seeker/BS-Seeker2/BS-Seeker3, Bismark, in mouse and

429  segemehl in human cancer cell lines. For the analysis of bisulfite sequence data, most of the

430  fundamentals of the chemical background and methylation principles are the same; however, the

431  major difference between the use of such tools in plants, and animals (specifically, in humans

432  and mouse) is the genome structure organization and the presence of predominantly more

433  CHG/CHH methylation contexts in plants. The most predominant context of DNA methylation in

434  mammals is the symmetric CG – estimated to be at ~70-80% of CG dinucleotides genome-wide

435  [140]. The mechanisms of regulation and function of DNA methylation vary in animals and plants

436  [141, 142]. These variations in regulation and function mechanism, coupled with genome

437  structure differences and complexity levels is a motivating factor for integrating small subtle

438  differences in mapping and analysis tools for epigenome data. Another important difference of

439  plants and animals is how they are able to demethylate their genome. So far, enzymes removing

440  directly the methyl group from cytosines have not been identified in plants, but they are important

441  components of mammalian DNA methylation homeostasis. Plants use either passive

442  mechanisms (not maintaining methylation during DNA replication) or base-excision and

443  subsequent repair for direct removal of methylated cytosines. Unlike with the human genome,

444  the CHG/CHH contexts which are more abundant in plants [143] need to be integrated into the

445  mapping and analysis of methylome data. Many plants have large and repetitive genomes

446  compared to that of humans. Such large genomes are a limiting factor in the analysis since they

447  require a lot of computational resources. The sequence mapping to references and statistical

448  computational time for large genomes such that of bread wheat (~17Gb) and barley (~5.3Gb) is

449  likely to scale linearly.

450  **Concluding remarks**

451 In the last decade, there has been tremendous progress in the development of tools for

452 analyzing epigenomic data; however, numerous challenges remain. For instance, the

453 visualization capacity of many tools remains either inadequate or lacks essential modules for

454 handling and displaying statistical outcomes from the resulting analysis. Additionally, the of these

455 tools to scale to handle large genomes remains an issue for further exploration. Technically,

456 most computational tools for analyzing epigenomic data perform well for datasets from

457 organisms with a genome size that is smaller than the human genome (~3Gb). For much larger

458 and complex genomes, more computational resources are required and the genome structure

459 (whether diploid, hexaploidy, or tetraploid) and repetitive nature of the genome has to be taken

460 into consideration during mapping to a reference genome. This is demonstrated in our example

461 where we compared the mapping efficiency for *Arabidopsis* and a wheat chromosome; however,

462 the complexity in genome structure, the presence of transposable elements, and the lack of

463 consistent gene annotations for some plants remain a major obstacle to advancing epigenetic

464 research.

465 In the next decade, there is likely to be a steady improvement in sequencing methods and

466 performance of already existing computational algorithms. Recently, it was shown that even well-

467 established sequencing methods might be prone to errors, leading to misleading results, e.g.

468 DNA immunoprecipitation sequencing (DIP-seq) [144]. Discovering and amending such errors

469 can lead to new findings from the previous studies and limit these errors' damage to future

470 studies. This will aid further epigenetic research not only in plants but also in life sciences in

471 general. Additionally, a few tools have the capability to effectively get more information out of

472 low-coverage data. Developing new tools or improving on existing ones to attain optimal results

473 using low coverage data and fewer replicates would save experiment and sequencing costs. A

474 high sequence coverage allows for good data quality and enables robust statistical analysis

475 [145]. Achieving high sequence coverage can be quite expensive and the minimum desired

476  coverage can depend on the research objectives at hand. Typically, a coverage value of 5-10X is

477  sufficient for many comparative studies and for achieving reliable methylation calls [145].

478  However, studies have demonstrated that coverage values as low as 2X is sufficient [146].

479  Accurate identification of DMRs in large samples, especially between multiple conditions,

480  remains a challenge – despite tremendous progress already made in this area.

**Acknowledgements**

**Conflicting interests**

485  No conflicting interests declared.

**Funding**

**Author Contributions**

489  J.O and T.N initiated the study, initiated the analysis framework, performed the data analysis,

490  simulations and drafted the manuscript. R.G stream-lined the manuscript content and write-up.

**Key points**

492  •  We introduce the concepts of epigenetics in plants and discuss commonly used tools –

493     with a focus on their capabilities.

494  •  Integration of bioinformatics tools needed to understand epigenomics datasets in crops.

495  •  The presence of repetitive elements in the genome influences the prediction of

496     methylated sites.

497 • We list the runtime and computational requirement for a small and large complex genome

498 and demonstrate their overlaps in four most applied tools.

499 • Different tools have different levels of asymmetry with regards to their mapping and

500 methylation call statistics.

501

## References

503 1. Costello, J.F. and C. Plass, *Methylation matters.* J Med Genet, 2001. **38**(5): p. 285-303.
504 2. Takuno, S., J.H. Ran, and B.S. Gaut, *Evolutionary patterns of genic DNA methylation*
505 *vary across land plants.* Nat Plants, 2016. **2**: p. 15222.
506 3. Bewick, A.J., et al., *On the origin and evolutionary consequences of gene body DNA*
507 *methylation.* Proc Natl Acad Sci U S A, 2016. **113**(32): p. 9111-6.
508 4. Bewick, A.J. and R.J. Schmitz, *Gene body DNA methylation in plants.* Curr Opin Plant
509 Biol, 2017. **36**: p. 103-110.
510 5. Wang, Y., et al., *Gene body methylation shows distinct patterns associated with different*
511 *gene origins and duplication modes and has a heterogeneous relationship with gene*
512 *expression in Oryza sativa (rice).* New Phytol, 2013. **198**(1): p. 274-83.
513 6. Bewick, A.J., et al., *Evolution of DNA Methylation across Insects.* Mol Biol Evol, 2017.
514 **34**(3): p. 654-665.
515 7. Lauss, K., et al., *Parental DNA Methylation States Are Associated with Heterosis in*
516 *Epigenetic Hybrids.* Plant Physiol, 2018. **176**(2): p. 1627-1645.
517 8. Kooke, R. and J.J. Keurentjes, *Epigenetic variation contributes to environmental*
518 *adaptation of Arabidopsis thaliana.* Plant Signal Behav, 2015. **10**(9): p. e1057368.
519 9. Cortijo, S., et al., *Mapping the epigenetic basis of complex traits.* Science, 2014.
520 **343**(6175): p. 1145-8.
521 10. Kakutani, T., *Epi-alleles in plants: inheritance of epigenetic information over generations.*
522 Plant Cell Physiol, 2002. **43**(10): p. 1106-11.
523 11. Quadrana, L. and V. Colot, *Plant Transgenerational Epigenetics.* Annu Rev Genet, 2016.
524 **50**: p. 467-491.
525 12. Hauser, M.T., et al., *Transgenerational epigenetic inheritance in plants.* Biochim Biophys
526 Acta, 2011. **1809**(8): p. 459-68.
527 13. Heard, E. and R.A. Martienssen, *Transgenerational epigenetic inheritance: myths and*
528 *mechanisms.* Cell, 2014. **157**(1): p. 95-109.
529 14. Grossniklaus, U., et al., *Transgenerational epigenetic inheritance: how important is it?*
530 Nat Rev Genet, 2013. **14**(3): p. 228-35.
531 15. Martienssen, R.A. and V. Colot, *DNA methylation and epigenetic inheritance in plants*
532 *and filamentous fungi.* Science, 2001. **293**(5532): p. 1070-4.
533 16. Deleris, A., et al., *Loss of the DNA methyltransferase MET1 Induces H3K9*
534 *hypermethylation at PcG target genes and redistribution of H3K27 trimethylation to*
535 *transposons in Arabidopsis thaliana.* PLoS Genet, 2012. **8**(11): p. e1003062.
536 17. Kim, Y.J., J. Lee, and K. Han, *Transposable Elements: No More 'Junk DNA'.* Genomics
537 Inform, 2012. **10**(4): p. 226-33.
538 18. Lister, R. and J.R. Ecker, *Finding the fifth base: genome-wide sequencing of cytosine*
539 *methylation.* Genome Res, 2009. **19**(6): p. 959-66.

19. Tomato Genome, C., *The tomato genome sequence provides insights into fleshy fruit evolution.* Nature, 2012. **485**(7400): p. 635-41.

20. Goff, S.A., et al., *A draft sequence of the rice genome (Oryza sativa L. ssp. japonica).* Science, 2002. **296**(5565): p. 92-100.

21. Lanciano, S. and M. Mirouze, *DNA Methylation in Rice and Relevance for Breeding.* Epigenomes, 2017. **1**(2): p. 10.

22. Cokus, S.J., et al., *Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA methylation patterning.* Nature, 2008. **452**(7184): p. 215-9.

23. Wicker, T., et al., *A unified classification system for eukaryotic transposable elements.* Nature Reviews Genetics, 2007. **8**: p. 973.

24. Tsukahara, S., et al., *Bursts of retrotransposition reproduced in Arabidopsis.* Nature, 2009. **461**: p. 423.

25. Eichten, S.R., et al., *DNA methylation profiles of diverse Brachypodium distachyon align with underlying genetic diversity.* Genome Res, 2016. **26**(11): p. 1520-1531.

26. Dubin, M.J., O. Mittelsten Scheid, and C. Becker, *Transposons: a blessing curse.* Curr Opin Plant Biol, 2018. **42**: p. 23-29.

27. Bourque, G., et al., *Ten things you should know about transposable elements.* Genome Biol, 2018. **19**(1): p. 199.

28. Hirsch, C.D. and N.M. Springer, *Transposable element influences on gene expression in plants.* Biochim Biophys Acta, 2017. **1860**(1): p. 157-165.

29. Dowen, R.H., et al., *Widespread dynamic DNA methylation in response to biotic stress.* Proc Natl Acad Sci U S A, 2012. **109**(32): p. E2183-91.

30. Dubin, M.J., et al., *DNA methylation in Arabidopsis has a genetic basis and shows evidence of local adaptation.* Elife, 2015. **4**: p. e05255.

31. Shen, X., et al., *Natural CMT2 variation is associated with genome-wide methylation changes and temperature seasonality.* PLoS Genet, 2014. **10**(12): p. e1004842.

32. Bossdorf O, P.D., Auge H, Schmid B., *Reduced competitive ability in an invasive plant.* Ecology Letters, 2004. **7**: p. 346–353.

33. Bouyer, D., et al., *DNA methylation dynamics during early plant life.* Genome Biol, 2017. **18**(1): p. 179.

34. Bartels, A., et al., *Dynamic DNA Methylation in Plant Growth and Development.* Int J Mol Sci, 2018. **19**(7).

35. Zhang, M., et al., *DNA cytosine methylation in plant development.* J Genet Genomics, 2010. **37**(1): p. 1-12.

36. Baurens, F.C., et al., *Genomic DNA methylation of juvenile and mature Acacia mangium micropropagated in vitro with reference to leaf morphology as a phase change marker.* Tree Physiol, 2004. **24**(4): p. 401-7.

37. Finnegan, E.J., W.J. Peacock, and E.S. Dennis, *Reduced DNA methylation in Arabidopsis thaliana results in abnormal plant development.* Proc Natl Acad Sci U S A, 1996. **93**(16): p. 8449-54.

38. Peng, H. and J. Zhang, *Plant genomic DNA methylation in response to stresses: potential applications and challenges in plant breeding.* Progress in Natural Science, 2009. **19**(9): p. 1037-1045.

39. Wang, W., et al., *Genome-wide differences in DNA methylation changes in two contrasting rice genotypes in response to drought conditions.* Frontiers in plant science, 2016. **7**: p. 1675.

40. Razin, A. and H. Cedar, *DNA methylation and gene expression.* Microbiol Rev, 1991. **55**(3): p. 451-8.

41. Zilberman, D., et al., *Genome-wide analysis of Arabidopsis thaliana DNA methylation uncovers an interdependence between methylation and transcription.* Nature genetics, 2007. **39**(1): p. 61.

591  42.  Li, X., et al., *Single-base resolution maps of cultivated and wild rice methylomes and regulatory roles of DNA methylation in plant gene expression.* BMC genomics, 2012. **13**(1): p. 300.

594  43.  Marx, V., *Genetics: profiling DNA methylation and beyond.* Nat Methods, 2016. **13**(2): p. 119-22.

596  44.  Eckhardt, F., et al., *DNA methylation profiling of human chromosomes 6, 20 and 22.* Nat Genet, 2006. **38**(12): p. 1378-85.

598  45.  Yong, W.S., F.M. Hsu, and P.Y. Chen, *Profiling genome-wide DNA methylation.* Epigenetics Chromatin, 2016. **9**: p. 26.

600  46.  Clark, S.J., et al., *High sensitivity mapping of methylated cytosines.* Nucleic Acids Res, 1994. **22**(15): p. 2990-7.

602  47.  Jeddeloh, J.A., J.M. Greally, and O.J. Rando, *Reduced-representation methylation mapping.* Genome Biol, 2008. **9**(8): p. 231.

604  48.  Schmidt, M., et al., *Plant-RRBS, a bisulfite and next-generation sequencing-based methylome profiling method enriching for coverage of cytosine positions.* BMC Plant Biol, 2017. **17**(1): p. 115.

607  49.  Patel, R.K. and M. Jain, *NGS QC Toolkit: a toolkit for quality control of next generation sequencing data.* PLoS One, 2012. **7**(2): p. e30619.

609  50.  Martin, M., *Cutadapt removes adapter sequences from high-throughput sequencing reads.* Bioinformatics in Action, 2012. **17**(1): p. 10-12.

611  51.  Bolger, A.M., M. Lohse, and B. Usadel, *Trimmomatic: a flexible trimmer for Illumina sequence data.* Bioinformatics, 2014. **30**(15): p. 2114-20.

613  52.  Merkel, A., et al., *gemBS - high throughput processing for DNA methylation data from Bisulfite Sequencing.* Bioinformatics, 2018.

615  53.  Marco-Sola, S., et al., *The GEM mapper: fast, accurate and versatile alignment by filtration.* Nat Methods, 2012. **9**(12): p. 1185-8.

617  54.  Liang, F., et al., *WBSA: web service for bisulfite sequencing data analysis.* PLoS One, 2014. **9**(1): p. e86707.

619  55.  Wreczycka, K., et al., *Strategies for analyzing bisulfite sequencing data.* J Biotechnol, 2017. **261**: p. 105-115.

621  56.  Mohn, F., et al., *Methylated DNA immunoprecipitation (MeDIP).* Methods Mol Biol, 2009. **507**: p. 55-64.

623  57.  Brinkman, A.B., et al., *Whole-genome DNA methylation profiling using MethylCap-seq.* Methods, 2010. **52**(3): p. 232-6.

625  58.  Booth, M.J., et al., *Oxidative bisulfite sequencing of 5-methylcytosine and 5-hydroxymethylcytosine.* Nat Protoc, 2013. **8**(10): p. 1841-51.

627  59.  Yu, M., et al., *Base-resolution analysis of 5-hydroxymethylcytosine in the mammalian genome.* Cell, 2012. **149**(6): p. 1368-80.

629  60.  Lu, X., et al., *Chemical modification-assisted bisulfite sequencing (CAB-Seq) for 5-carboxylcytosine detection in DNA.* J Am Chem Soc, 2013. **135**(25): p. 9315-7.

631  61.  Wang, X.L., et al., *Genome-wide mapping of 5-hydroxymethylcytosine in three rice cultivars reveals its preferential localization in transcriptionally silent transposable element genes.* J Exp Bot, 2015. **66**(21): p. 6651-63.

634  62.  Shi, D.Q., et al., *New Insights into 5hmC DNA Modification: Generation, Distribution and Function.* Front Genet, 2017. **8**: p. 100.

636  63.  Erdmann, R.M., et al., *5-hydroxymethylcytosine is not present in appreciable quantities in Arabidopsis DNA.* G3 (Bethesda), 2014. **5**(1): p. 1-8.

638  64.  Krueger, F. and S.R. Andrews, *Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications.* Bioinformatics, 2011. **27**(11): p. 1571-2.

640  65.  Xi, Y. and W. Li, *BSMAP: whole genome bisulfite sequence MAPping program.* BMC Bioinformatics, 2009. **10**: p. 232.

66. Langmead, B., et al., *Ultrafast and memory-efficient alignment of short DNA sequences to the human genome.* Genome Biol, 2009. **10**(3): p. R25.

67. Hansen, K.D., B. Langmead, and R.A. Irizarry, *BSmooth: from whole genome bisulfite sequencing reads to differentially methylated regions.* Genome Biol, 2012. **13**(10): p. R83.

68. Langmead, B. and S.L. Salzberg, *Fast gapped-read alignment with Bowtie 2.* Nat Methods, 2012. **9**(4): p. 357-9.

69. Chen, P.Y., S.J. Cokus, and M. Pellegrini, *BS Seeker: precise mapping for bisulfite sequencing.* BMC Bioinformatics, 2010. **11**: p. 203.

70. Guo, W., et al., *BS-Seeker2: a versatile aligning pipeline for bisulfite sequencing data.* BMC Genomics, 2013. **14**: p. 774.

71. Huang, K.Y.Y., Y.J. Huang, and P.Y. Chen, *BS-Seeker3: ultrafast pipeline for bisulfite sequencing.* BMC Bioinformatics, 2018. **19**(1): p. 111.

72. Harris, E.Y., R. Ounit, and S. Lonardi, *BRAT-nova: fast and accurate mapping of bisulfite-treated reads.* Bioinformatics, 2016. **32**(17): p. 2696-8.

73. Chen, H., A.D. Smith, and T. Chen, *WALT: fast and accurate read mapping for bisulfite sequencing.* Bioinformatics, 2016. **32**(22): p. 3507-3509.

74. Adusumalli, S., et al., *Methodological aspects of whole-genome bisulfite sequencing analysis.* Brief Bioinform, 2015. **16**(3): p. 369-79.

75. Shafi, A., et al., *A survey of the approaches for identifying differential methylation using bisulfite sequencing data.* Brief Bioinform, 2017.

76. Zhang, L., et al., *A nonparametric Bayesian approach for clustering bisulfate-based DNA methylation profiles.* BMC Genomics, 2012. **13 Suppl 6**: p. S20.

77. Cedoz, P.L., et al., *MethylMix 2.0: an R package for identifying DNA methylation genes.* Bioinformatics, 2018.

78. Widman, N., et al., *Epigenetic differences between shoots and roots in Arabidopsis reveals tissue-specific regulation.* Epigenetics, 2014. **9**(2): p. 236-42.

79. Turco, G.M., et al., *DNA methylation and gene expression regulation associated with vascularization in Sorghum bicolor.* New Phytol, 2017. **214**(3): p. 1213-1229.

80. Gardiner, L.J., et al., *Hidden variation in polyploid wheat drives local adaptation.* Genome Res, 2018. **28**(9): p. 1319-1332.

81. Lea, A.J., J. Tung, and X. Zhou, *A Flexible, Efficient Binomial Mixed Model for Identifying Differential DNA Methylation in Bisulfite Sequencing Data.* PLoS Genet, 2015. **11**(11): p. e1005650.

82. Li, S., et al., *An optimized algorithm for detecting and annotating regional differential methylation.* BMC Bioinformatics, 2013. **14 Suppl 5**: p. S10.

83. Akalin, A., et al., *methylKit: a comprehensive R package for the analysis of genome-wide DNA methylation profiles.* Genome Biol, 2012. **13**(10): p. R87.

84. Tian, Y., et al., *ChAMP: updated methylation analysis pipeline for Illumina BeadChips.* Bioinformatics, 2017. **33**(24): p. 3982-3984.

85. Jaffe, A.E., et al., *Bump hunting to identify differentially methylated regions in epigenetic epidemiology studies.* Int J Epidemiol, 2012. **41**(1): p. 200-9.

86. Butcher, L.M. and S. Beck, *Probe Lasso: a novel method to rope in differentially methylated regions with 450K DNA methylation data.* Methods, 2015. **72**: p. 21-8.

87. Ward, J.H., *Hierarchical Grouping to Optimize an Objective Function.* Journal of the American Statistical Association, 1963. **58**(301): p. 236-244.

88. J. van der Laan, M. and K.S. Pollard, *A new algorithm for hybrid hierarchical clustering with visualization and the bootstrap.* Journal of Statistical Planning and Inference, 2003. **117**(2): p. 275-303.

89. Gardiner, L.-J., et al., *A genome-wide survey of DNA methylation in hexaploid wheat.* Genome biology, 2015. **16**(1): p. 273.

90. Zhang, Y., et al., *Large-scale comparative epigenomics reveals hierarchical regulation of non-CG methylation in Arabidopsis.* Proc Natl Acad Sci U S A, 2018. **115**(5): p. E1069-E1074.

91. Mager, S. and U. Ludewig, *Massive Loss of DNA Methylation in Nitrogen-, but Not in Phosphorus-Deficient Zea mays Roots Is Poorly Correlated With Gene Expression Differences.* Front Plant Sci, 2018. **9**: p. 497.

92. Lindroth, A.M., et al., *Requirement of CHROMOMETHYLASE3 for maintenance of CpXpG methylation.* Science, 2001. **292**(5524): p. 2077-80.

93. Law, J.A. and S.E. Jacobsen, *Establishing, maintaining and modifying DNA methylation patterns in plants and animals.* Nat Rev Genet, 2010. **11**(3): p. 204-20.

94. Barrero, M.J., S. Boue, and J.C. Izpisua Belmonte, *Epigenetic mechanisms that regulate cell identity.* Cell Stem Cell, 2010. **7**(5): p. 565-70.

95. Juhling, F., et al., *metilene: fast and sensitive calling of differentially methylated regions from bisulfite sequencing data.* Genome Res, 2016. **26**(2): p. 256-62.

96. Catoni, M., et al., *DMRcaller: a versatile R/Bioconductor package for detection and visualization of differentially methylated regions in CpG and non-CpG contexts.* Nucleic Acids Res, 2018. **46**(19): p. e114.

97. Hebestreit, K., M. Dugas, and H.U. Klein, *Detection of significantly differentially methylated regions in targeted bisulfite sequencing data.* Bioinformatics, 2013. **29**(13): p. 1647-53.

98. Kurdyukov, S. and M. Bullock, *DNA Methylation Analysis: Choosing the Right Method.* Biology (Basel), 2016. **5**(1).

99. Robinson, M.D., et al., *Statistical methods for detecting differentially methylated loci and regions.* Front Genet, 2014. **5**: p. 324.

100. Deaton, A.M. and A. Bird, *CpG islands and the regulation of transcription.* Genes Dev, 2011. **25**(10): p. 1010-22.

101. Ashikawa, I., *Gene-associated CpG islands in plants as revealed by analyses of genomic sequences.* Plant J, 2001. **26**(6): p. 617-25.

102. Ashikawa, I., *Gene-associated CpG islands and the expression pattern of genes in rice.* DNA Res, 2002. **9**(4): p. 131-4.

103. Bock, C., et al., *BiQ Analyzer: visualization and quality control for DNA methylation data from bisulfite sequencing.* Bioinformatics, 2005. **21**(21): p. 4067-8.

104. Lin, X., et al., *BSeQC: quality control of bisulfite sequencing experiments.* Bioinformatics, 2013. **29**(24): p. 3227-9.

105. Lukauskas, S., et al., *DGW: an exploratory data analysis tool for clustering and visualisation of epigenomic marks.* BMC Bioinformatics, 2016. **17**(Suppl 16): p. 447.

106. Chari, R., et al., *Integrating the multiple dimensions of genomic and epigenomic landscapes of cancer.* Cancer Metastasis Rev, 2010. **29**(1): p. 73-93.

107. Sunami, E., et al., *LINE-1 hypomethylation during primary colon cancer progression.* PLoS One, 2011. **6**(4): p. e18884.

108. Ji, H., et al., *APOE hypermethylation is significantly associated with coronary heart disease in males.* Gene, 2018. **689**: p. 84-89.

109. Ghavifekr Fakhr, M., et al., *DNA methylation pattern as important epigenetic criterion in cancer.* Genet Res Int, 2013. **2013**: p. 317569.

110. Finnegan, E.J., et al., *DNA methylation and the promotion of flowering by vernalization.* Proc Natl Acad Sci U S A, 1998. **95**(10): p. 5824-9.

111. Eichten, S.R. and N.M. Springer, *Minimal evidence for consistent changes in maize DNA methylation patterns following environmental stress.* Front Plant Sci, 2015. **6**: p. 308.

112. Li, J., et al., *Global DNA methylation variations after short-term heat shock treatment in cultured microspores of Brassica napus cv. Topas.* Sci Rep, 2016. **6**: p. 38401.

743  113.  Sun, S., A. Noviski, and X. Yu, *MethyQA: a pipeline for bisulfite-treated methylation*
744        *sequencing quality assessment.* BMC Bioinformatics, 2013. **14**: p. 259.

745  114.  Worm Orntoft, M.B., et al., *Comparative analysis of 12 different kits for bisulfite*
746        *conversion of circulating cell-free DNA.* Epigenetics, 2017. **12**(8): p. 626-636.

747  115.  Holmes, E.E., et al., *Performance evaluation of kits for bisulfite-conversion of DNA from*
748        *tissues, cell lines, FFPE tissues, aspirates, lavages, effusions, plasma, serum, and urine.*
749        PLoS One, 2014. **9**(4): p. e93933.

750  116.  Liu, Y.Y. and H.M. Cui, *The method of estimating bisulfite conversion rate in DNA*
751        *methylation analysis.* Yi Chuan, 2015. **37**(9): p. 939-44.

752  117.  International Wheat Genome Sequencing, C., et al., *Shifting the limits in wheat research*
753        *and breeding using a fully annotated reference genome.* Science, 2018. **361**(6403).

754  118.  Kawakatsu, T., et al., *Epigenomic Diversity in a Global Collection of Arabidopsis thaliana*
755        *Accessions.* Cell, 2016. **166**(2): p. 492-505.

756  119.  Schmitz, R.J., et al., *Patterns of population epigenomic diversity.* Nature, 2013.
757        **495**(7440): p. 193-8.

758  120.  Hoffmann, S., et al., *Fast Mapping of Short Sequences with Mismatches, Insertions and*
759        *Deletions Using Index Structures.* PLOS Computational Biology, 2009. **5**(9): p. e1000502.

760  121.  Tran, H., et al., *Objective and comprehensive evaluation of bisulfite short read mapping*
761        *tools.* Adv Bioinformatics, 2014. **2014**: p. 472045.

762  122.  Lee, J.H., S.J. Park, and N. Kenta, *An integrative approach for efficient analysis of whole*
763        *genome bisulfite sequencing data.* BMC Genomics, 2015. **16 Suppl 12**: p. S14.

764  123.  Hatem, A., et al., *Benchmarking short sequence mapping tools.* BMC Bioinformatics,
765        2013. **14**: p. 184.

766  124.  Harris, E.Y., et al., *BRAT: bisulfite-treated reads analysis tool.* Bioinformatics, 2010.
767        **26**(4): p. 572-3.

768  125.  Dinh, H.Q., et al., *Advanced methylome analysis after bisulfite deep sequencing: an*
769        *example in Arabidopsis.* PLoS One, 2012. **7**(7): p. e41528.

770  126.  Chinnusamy, V. and J.K. Zhu, *Epigenetic regulation of stress responses in plants.* Curr
771        Opin Plant Biol, 2009. **12**(2): p. 133-9.

772  127.  Kumar, S., *Epigenomics of Plant Responses to Environmental Stress.* Epigenomes,
773        2018. **2**(1): p. 6.

774  128.  Lamke, J. and I. Baurle, *Epigenetic and chromatin-based mechanisms in environmental*
775        *stress adaptation and stress memory in plants.* Genome Biol, 2017. **18**(1): p. 124.

776  129.  Gutzat, R. and O. Mittelsten Scheid, *Epigenetic responses to stress: triple defense?* Curr
777        Opin Plant Biol, 2012. **15**(5): p. 568-73.

778  130.  Kinoshita, T. and M. Seki, *Epigenetic memory for stress response and adaptation in*
779        *plants.* Plant Cell Physiol, 2014. **55**(11): p. 1859-63.

780  131.  Yang, Y. and Y. Guo, *Unraveling salt stress signaling in plants.* J Integr Plant Biol, 2018.

781  132.  Abhinandan, K., et al., *Abiotic Stress Signaling in Wheat - An Inclusive Overview of*
782        *Hormonal Interactions During Abiotic Stress Responses in Wheat.* Front Plant Sci, 2018.
783        **9**: p. 734.

784  133.  Gallusci, P., et al., *DNA Methylation and Chromatin Regulation during Fleshy Fruit*
785        *Development and Ripening.* Front Plant Sci, 2016. **7**: p. 807.

786  134.  Manning, K., et al., *A naturally occurring epigenetic mutation in a gene encoding an SBP-*
787        *box transcription factor inhibits tomato fruit ripening.* Nat Genet, 2006. **38**(8): p. 948-52.

788  135.  Omidvar, V. and M. Fellner, *DNA methylation and transcriptomic changes in response to*
789        *different lights and stresses in 7B-1 male-sterile tomato.* PLoS One, 2015. **10**(4): p.
790        e0121864.

791  136.  Farinati, S., et al., *Rosaceae Fruit Development, Ripening and Post-harvest: An*
792        *Epigenetic Perspective.* Front Plant Sci, 2017. **8**: p. 1247.

137.    Boyko, A. and I. Kovalchuk, *Epigenetic control of plant stress response.* Environ Mol Mutagen, 2008. **49**(1): p. 61-72.

138.    White, N.R. and R.J. Barfield, *Playback of female rat ultrasonic vocalizations during sexual behavior.* Physiol Behav, 1989. **45**(2): p. 229-33.

139.    Xu, S. and K. Chong, *Remembering winter through vernalisation.* Nat Plants, 2018. **4**(12): p. 997-1009.

140.    Ehrlich, M., et al., *Amount and distribution of 5-methylcytosine in human DNA from different types of tissues of cells.* Nucleic Acids Res, 1982. **10**(8): p. 2709-21.

141.    He, X.J., T. Chen, and J.K. Zhu, *Regulation and function of DNA methylation in plants and animals.* Cell Res, 2011. **21**(3): p. 442-65.

142.    Yi, S.V., *Insights into Epigenome Evolution from Animal and Plant Methylomes.* Genome Biol Evol, 2017. **9**(11): p. 3189-3201.

143.    Su, Z., L. Han, and Z. Zhao, *Conservation and divergence of DNA methylation in eukaryotes: new insights from single base-resolution DNA methylomes.* Epigenetics, 2011. **6**(2): p. 134-140.

144.    Lentini, A., et al., *A reassessment of DNA-immunoprecipitation-based genomic profiling.* Nat Methods, 2018.

145.    Ziller, M.J., et al., *Coverage recommendations for methylation analysis by whole-genome bisulfite sequencing.* Nat Methods, 2015. **12**(3): p. 230-2, 1 p following 232.

146.    Li, Q., et al., *Examining the Causes and Consequences of Context-Specific Differential DNA Methylation in Maize.* Plant Physiol, 2015. **168**(4): p. 1262-74.

816 **Figures**

817 **Figure 1. Selection of epigenomics tools**.

818 Figure panels A and B: Results of the calculation user times for four common tools, Bismark,

819 BSMap, BS-Seeker3, and segemehl. We used data for *Arabidopsis thaliana* and chromosome

820 1A in bread wheat (*Triticum aestivum*). *n.a*: values not available. Figure panels C and D: Overlap

821 of detected sites in the two reference genomes for the four mapping tools.

822

823 **Figure 2. Precision and sensitivity analysis**.

824 Precision and sensitivity analysis for the *A. thaliana* data based on read mapping of simulated

825 reads using the tool by Sherman (https://www.bioinformatics.babraham.ac.uk/projects/sherman/)

826 – with the parameters (CG=24, CH=8, e=0.5). (A) There is a large difference in the sensitivity of

827 the four tools. BS-Seeker3 was the least sensitive (sensitivity averaging ~ 48%) – Bismark was

828 the most sensitive (sensitivity ~99.9%). The sensitivity values for BSMap and segemehl

829 averaged ~97% and 90%, respectively. (B) For bread wheat (*T. aestivum*), BSMap appears to be

830 marginally less precise and less sensitive than segemehl. There is consistency in the precision

831 and sensitivity values for the subgenomes A, B and D in chromosome 1 of *T. aestivum*. Overall,

832 the results from both (A) and (B) are in agreement. Notably, BS-Seeker3 has a wide range of

833 precision compared to the other three tools. Each data point represents the precision-sensitivity

834 value based on a simulation run for an individual tool. The precision and sensitivity values for

835 Bismark, BSMap, BS-Seeker3 and segemehl averaged approximately (99%, 99%), (94%, 82%),

836 (86%, 38%) and (97%, 87%), respectively. Five (5) simulation runs were performed for each tool

837 – one for each of the *A. thaliana* chromosomes. The elliptical rings around each set of (same

838 colored) data points represent the confidence bounds.

839

840   **Figure 3. Memory footprint analysis for the four tools – benchmarked on five genomes.**

841   (A) Barplots showing variation in attained memory footprint between the tools benchmarked on

842   different genomes. (B to E) Correlation analysis of genome size and memory footprint analysis. A

843   benchmark of the four tools, (B) BSMap, (C) BS-Seeker3, (D) Bismark, and (E) segemehl. The

844   genome sizes are all significantly correlated to the memory footprint analysis (*p*-values < 0.05).

845   Red dotted line: fitted regression line, green-dots: data points.

846

847   **Table 1.** Examples of some down-stream analysis software.

848

849   **Supporting information**

850   **Table S1. A selection of popular packages and tools for epigenome data analysis.**

851   Unranked list compiled based on high (≥100) citation and usage (October 2018). Most of these

852   tools are freely available for download (non-commercial) and some are embedded into a web-

853   server.

854

855   **Table S2. Summary of key attributes and parameters for the four benchmarked tools.**

856

**(A)** *Arabidopsis thaliana* dataset

| Amount of threads | BSMap | BS-Seeker3 | Bismark | segemehl | |
|---|---|---|---|---|---|
| 1 | 1m 32s<br>1m 33s | 1m 8s<br>8m 12s | 10m 10s<br>29m 56s | 185m 38s<br>181m 17s | real time<br>user time |
| 5 | 0m 25s<br>1m 34s | n.a | 3m 51s<br>33m 52s | 44m 9s<br>208m 31s | |
| 10 | 0m 17s<br>1m 36s | n.a | 3m 35s<br>34m 55s | 24m 15s<br>207m 30s | |
| 15 | 0m 15s<br>1m 38s | n.a | 3m 36s<br>35m 47s | 17m 33s<br>207m 58s | |
| 20 | 0m 15s<br>1m 39s | n.a | 3m 35s<br>36m 7s | 17m 6s<br>207m 14s | |

**(B)** *Triticum aestivum* Chr 1A

| Amount of threads | BSMap | BS-Seeker3 | Bismark | segemehl | |
|---|---|---|---|---|---|
| 1 | 72m 44s<br>72m 40s | 4m 14s<br>23m 35s | 377m 20s<br>776m 6s | 1778m 3s<br>1711m 53s | real time<br>user time |
| 5 | 15m 35s<br>75m 36s | n.a | 85m 55s<br>824m 40s | 466m 24s<br>1990m 40s | |
| 10 | 8m 12s<br>76m 53s | n.a | 56m 41s<br>837m 42s | 280m 57s<br>1995m 15s | |
| 15 | 5m 58s<br>79m 49s | n.a | 57m 19s<br>850m 17s | 224m 26s<br>2026m 37s | |
| 20 | 5m 31s<br>79m 37s | n.a | 58m 21s<br>861m 8s | 214m 38s<br>2026m 9s | |



Selection of epigenomics tools. Figure panels A and B: Results of the calculation user times for four common tools, Bismark, BSMap, BS-Seeker3, and segemehl. We used data for Arabidopsis thaliana and chromosome 1A in bread wheat (Triticum aestivum). n.a: values not available. Figure panels C and D: Overlap of detected sites in the two reference genomes for the four mapping tools.

126x95mm (300 x 300 DPI)

Precision and sensitivity analysis. Precision and sensitivity analysis for the A. thaliana data based on read mapping of simulated rea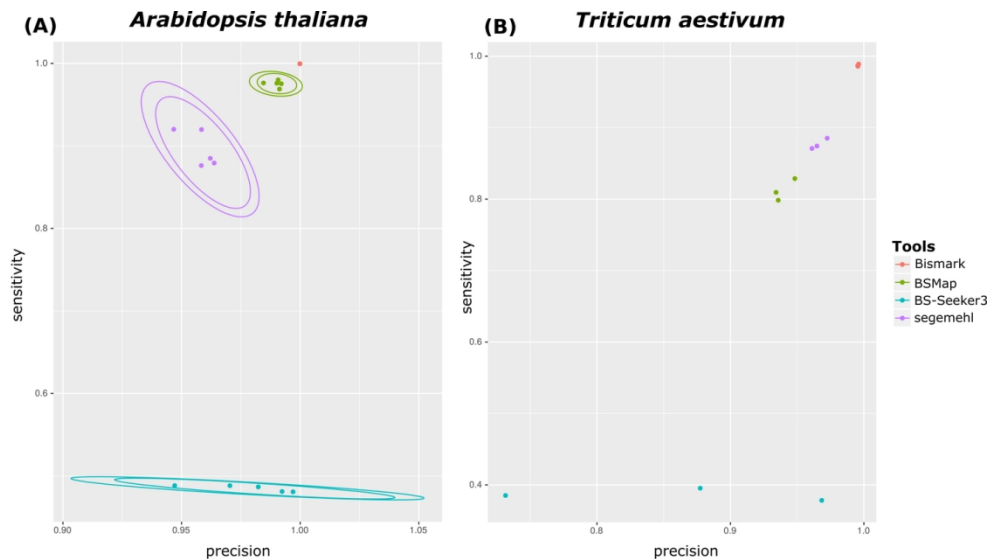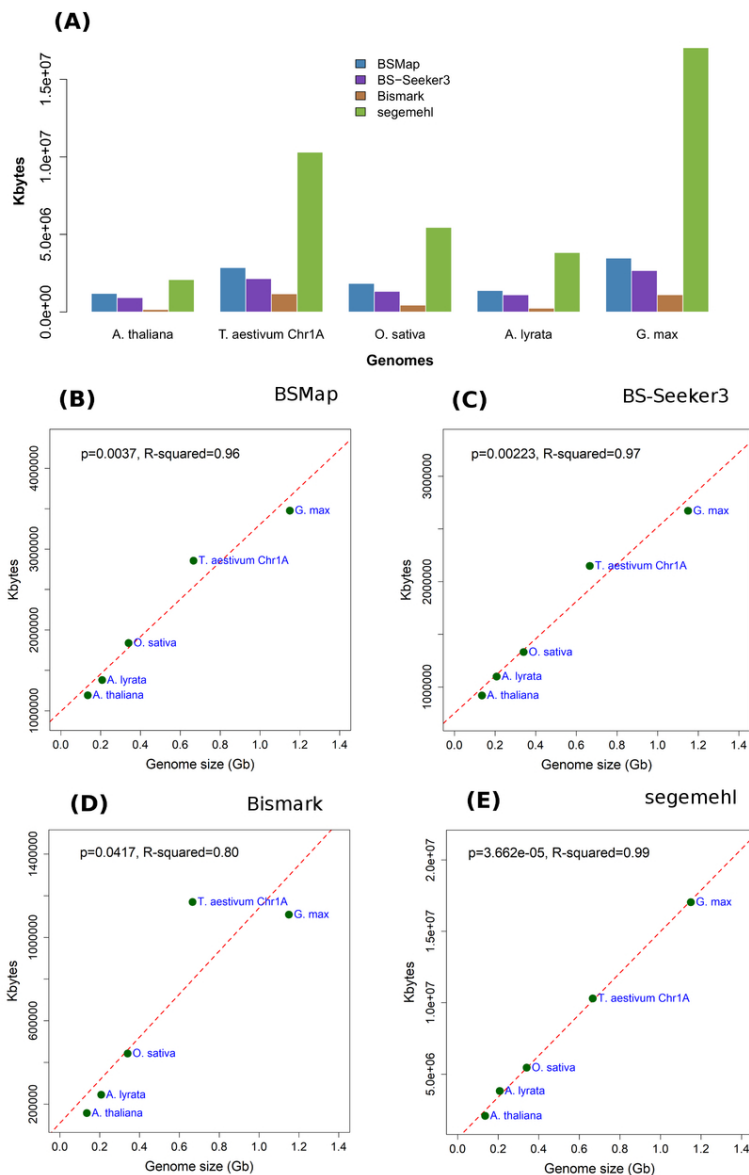ds using the tool by Sherman (https://www.bioinformatics.babraham.ac.uk/projects/sherman/) – with the parameters (CG=24, CH=8, e=0.5). (A) There is a large difference in the sensitivity of the four tools. BS-Seeker3 was the least sensitive (sensitivity averaging ~ 48%) – Bismark was the most sensitive (sensitivity ~99.9%). The sensitivity values for BSMap and segemehl averaged ~97% and 90%, respectively. (B) For bread wheat (T. aestivum), BSMap appears to be marginally less precise and less sensitive than segemehl. There is consistency in the precision and sensitivity values for the subgenomes A, B and D in chromosome 1 of T. aestivum. Overall, the results from both (A) and (B) are in agreement. Notably, BS-Seeker3 has a wide range of precision compared to the other three tools. Each data point represents the precision-sensitivity value based on a simulation run for an individual tool. The precision and sensitivity values for Bismark, BSMap, BS-Seeker3 and segemehl averaged approximately (99%, 99%), (94%, 82%), (86%, 38%) and (97%, 87%), respectively. Five (5) simulation runs were performed for each tool – one for each of the A. thaliana chromosomes. The elliptical rings around each set of (same colored) data points represent the confidence bounds.

127x70mm (300 x 300 DPI)

Memory footprint analysis for the four tools – benchmarked on five genomes.
(A) Barplots showing variation in attained memory footprint between the tools benchmarked on different genomes. (B to E) Correlation analysis of genome size and memory footprint analysis. A benchmark of the four tools, (B) BSMap, (C) BS-Seeker3, (D) Bismark, and (E) segemehl. The genome sizes are all significantly correlated to the memory footprint analysis (p-values < 0.05). Red dotted line: fitted regression line, green-dots: data points.

71x107mm (300 x 300 DPI)

| Tool | Citation and descriptions |
|---|---|
| **ADMIRE**: Analysis and visualization of differential methylation in genomic regions using the Infinium HumanMethylation450 Assay | Preussner *et al*. [109]; Online and offline. Adds experimental settings, quality control, automatic filtering, normalization, multiple testing, and differential analyses genome-browser tracks, table outputs, summary files. |
| **BATMAN**: Bayesian automated metabolite analyser for NMR spectra | Hao *et al*. [110]; Uses Markov chain Monte Carlo algorithm for sampling. Bayesian based approach. |
| **KEGG**: Gene Ontology Pathways | It is a database for mining and analysis of high-level functions. KEGG enables analysis and data mining on different biological scales (e.g. cellular and molecular-level information, whole organism, at ecosystem level, etc – using data from high-throughput experiments; see https://www.genome.jp/kegg/). |
| **IPA**: Ingenuity Pathway Analysis | Krämer *et al*. [111]; Platform enables exploration and visualization of complex omics data (e.g. microarrays including miRNA, metabolomics, proteomics, RNA-seq, small RNA-seq and SNP, and small scale experiments). See https://www.qiagenbioinformatics.com/ |
| **DAVID**: Database for Annotation, Visualization and Integrated Discovery | Huang *et al*. [112]; DAVID enables pathway mining and gene function classification. Input is gene list from high-throughput genomic experiments; https://david.ncifcrf.gov/ |

| Category | Tool | Year | Software | Mapping tool | Method | Input | Output | Reference/web-page |
|---|---|---|---|---|---|---|---|---|
| Mapping | Bismark | 2011 | Perl | Bowtie | Integrates Bowtie, running four alignment processes simultaneously | FASTA/FASTQ | SAM format, BAM format, one entry (or line) per cytosine | https://github.com/FelixKrueger/Bismark |
| | BSmooth | 2012 | R/Perl | Merman, Bowtie, | Identification of DMRs, Walch t-test, improving previous work based on Fisher's exact test to simulate biological replicates using merman mapper | FASTQ | DMRs | https://github.com/BenLangmead/bsmooth-align |
| | | | | | | | | https://github.com/hansenlab/bsseq |
| | BSMap | 2009 | C++ | SOAP | Using HASH table seeding + Bitwise masking, bisulfite seq. data mapping program | FASTA/FASTQ/BAM - supports paired end reads | BSP/SAM/BAM | https://code.google.com/archive/p/bsmap/ |
| | RRBSMap | 2012 | C++ | RRBSMAP | RRBS short read alignment tool | BAM files | n.a. | http://rrbsmap.computational-epigenetics.org/ |
| | BS Seeker | 2010 | Python | Bowtie | Pipeline for mapping bisulfite sequence data, genome indexing. Accepts both RRB | FASTA/FASTQ, qseq, pure sequence, IGV input | BAM, SAM, BS_seeker and WIG files | https://guoweilong.github.io/BS_Seeker2/index.html |
| | BS Seeker2 | 2013 | Python | Bowtie2 | Pipeline for mapping bisulfite sequence data, genome indexing. Accepts both RRB | FASTA/FASTQ, qseq, pure sequence, IGV input | BAM, SAM, BS_seeker and WIG files | https://guoweilong.github.io/BS_Seeker2/index.html |
| | BS Seeker3 | 2018 | Python | SNAP-aligner | Pipeline for mapping bisulfite sequence data, genome indexing. Accepts both RRB | FASTA/FASTQ, qseq, pure sequence, IGV input | BAM, SAM, BS_seeker and WIG files | https://github.com/khuang28jhu/bs3 |
| | Metilene | 2015 | Perl | n.a. | DMR finder (from whole genome and targeted sequencing data) | FASTQ | DMRs | http://genome.cshlp.org/content/26/2/256.short |
| | segemehl | 2012 | Perl | segemehl | read mapper. Analysis of: COV, MET, TXN, SNP and CNV. | FASTQ | SAM format | http://www.bioinf.uni-leipzig.de/Software/segemehl/ |
| | BRAT-BW | 2012 | Free and Open Sour | FM-index (Burrow | Mapping of bisulfite reads, supports paired-end libraries, indels, mismatches, | FASTQ for reads, FASTA for reference sequence | Text files of mapping results | http://compbio.cs.ucr.edu/brat/ |
| | BRAT-nova | 2016 | Free and Open Sour | FM-index (Burrow | Mapping of bisulfite reads, improved implementation of the mapping tool BRAT-BW | FASTQ for reads, FASTA for reference sequence | Text files of mapping results | http://compbio.cs.ucr.edu/brat/ |
| | WALT | 2016 | Free and Open Source under GPLv3 | | Mapping bisulfite sequencing reads | FASTQ for reads, FASTA for reference sequence | SAM or MR files | https://github.com/smithlabcode/walt |
| **Statistical analysis** | | | | | | | | |
| | MethGo | 2015 | Python | n.a | global and gene level scale methylation pattern around TSS sites. Accepts both RR | FASTA, BAM, GTF and CGmap | SNP, CNV tables, methylation profile summaries/plots/tal | http://paoyangchen-laboratory.github.io/methgo/ |
| | EpiGRAPH | 2009 | Java | n.a | software for genome and epigenome analysis, uses machine learning algorithms | FASTA/FASTQ, genomic seq data | Enable prediction of genomic regions having similar char | https://epigraph.mpi-inf.mpg.de/WebGRAPH/ |
| | CyMATE | 2008 | Perl and C | n.a. | Unique mapping tool for CpG and non-CpG methylation (web-based tool) | multiple sequence alignment | Text files of per sequence mC, per position mC, global m | http://www.cymate.org/ |
| | MethylMapper | 2005 | Perl | MethylMapper | High through-put mapping (web-based tool), performs QC analysis | DNA methylation seq. data | File with counts/tallies of methylated sites | http://methylmapper.wiki/Home/ |
| | RnBeads | 2014 | R | n.a., using bed file | Supported assays: WGBS, RRBS. Data QC and filtering, DMR finder, Data explor | Illumina microarray platform bisulfite sequencing | bed, bigBed and bigWig file | https://rnbeads.org/ |
| | BISMA | 2010 | PHP code, Perl and | Uses ClustalW alg | Analysis of bisulfite sequence data. Supports analysis of repetitive sequences (we | ABI, text and single multi-FASTA file formats | Outputs multiple sequence alignment. Web-presented ou | http://services.ibc.uni-stuttgart.de/BDPC/BISMA/ |
| | BSPAT | 2015 | Java, Tomcat | Bowtie | Online service to analyze methylation patterns in bisulfite sequencing data. Has in | FASTA, FASTQ | SAM | https://github.com/lancelothk/BSPAT |
| | MethylMix | 2015 | R | n.a | Detects hyper- and hypomethylated regions. Uses Bayesian Information Criterion (BIC) to select number of methylation states by iter | Differentially methylated regions (DMRs) | https://www.bioconductor.org/packages/release/bioc/html/MethylMix.html |
| | bumphunter | 2012 | R | n.a. | Detection of DMRs | FASTA, FASTQ | BED, BigWig files | https://github.com/rafalab/bumphunter |
| | DMAP | 2014 | Bismark alignment | n.a | Differentially methylated region detection. Accepts both RRBSand WGBS data | SAM file from aligner tool | DMRs, identifies genes and CpG features; distances to D | http://biochem.otago.ac.nz/research/databases-software/ |
| **Complete pipeline** | | | | | | | | |
| | SAAP-RRBS | 2012 | BSMap: modules de | uses hashing/bitw | Includes FASTQC, duplicate read removal, read alignment, and methylation calls | FASTQ | Bed files with CpG sites, annotation files | http://bioinformaticstools.mayo.edu/research/saap-rrbs/ |
| | gemBS | 2018 | C, Python | n.a. | Analysis of whole genome bisulfite sequence data (both WGBS, RRBS) | FASTQ/FASTA, SAM/BAM files | Supports DMRs, Yes (bigWig, bedGraph) | http://statgen.cnag.cat/GEMBS/ |
| | BiQ Analyzer HiMod | 2014 | web server | n.a. | Upgrade of BiQ Analyzer HT | FASTA, FASTQ | BED, bigBed, ... and BigWig files | https://biq-analyzer-himod.bioinf.mpi-inf.mpg.de/ |
| | BiQ Analyzer HT | 2011 | web server | n.a. | Bisulfite sequence quality assessement tool. Analysis of unique sequences | FASTA or BAM files | Outputs multiple sequence alignment. Web-presented ou | https://biq-analyzer-ht.bioinf.mpi-inf.mpg.de/ |
| | QUMA | 2008 | n.a. | EMBOSS package | Interactive web-based tool | FASTQ data, bisulfite sequence | Outputs multiple sequence alignment, statistics summary | http://quma.cdb.riken.jp/ |

| Abbreviations | Explanation |
|---|---|
| RRBS | Reduced representation bisulfite sequencing |
| WGBS | Whole-genome bisulfite sequence |
| SAM | Sequence Alignment Map |
| BAM | Binary Alignment Map |
| TFBS | Transcription factor binding site |
| COV | Coverage distribution of methylation sites |
| MET | Methylation profiling |
| TXN | Cytosine methylation levels at transcription factor binding sites (TFBSs) |
| SNP | Single-nucleotide polymorphism |
| CNV | Copy number variation |
| GTF | Gene transfer format |
| HT | High through-put |
| QC | Quality control |
| DMRs | Differentially methylated regions |
| SAAP-RRBS | Streamlined Analysis and Annotation Pipeline for RRBS data |
| n.a | Not available (no explicitely specified) |

**Comparison of fearures between the four tools**

| Feature | BSMap | BS-Seeker3 | Bismark | segemehl |
|---|---|---|---|---|
| Allows for multiple threads | Yes | No | Yes | Yes |
| Supports single-end(SE)/paired-end(PE) reads | Yes/Yes | Yes/Yes | Yes/Yes | Yes/Yes |
| Variable read length (SE/PE) | Yes/Yes | Yes/Yes | Yes/Yes | Yes/Yes |
| Allows for mismatches during mapping | Yes | Yes | Yes | Yes |
| Allows for adaptor trimming | Yes | No | No | No |
| Supports gapped alignments | Yes | Yes | Yes | Yes |
| Supports RRBS/WGBS | Yes/Yes | Yes/Yes | Yes/Yes | Yes |
| Outputs methylation by context (CpG/CHG/CHH) | Yes | Yes | Yes | Yes |
| Multiple adjustable mapping parameters (e.g. seed size, byte size, ...) | Yes | Yes | Yes | Yes |
| Strategy used for mapping | wild-card | 3-letter | 3-letter | 3-letter |
| Supports directional/non-directional libraries | Yes/Yes | Yes/Yes | Yes/Yes | Yes/Yes |
| Allows for priliminary quality control analysis | not specified | Yes | Yes | not specified |
| Provides tabular/visual summary mapping statistics | Yes | Yes | Yes | Yes |
| Maximum read length allowed | 144nt | not specified | variable size | not specified |

| Abbreviations | Explanation |
|---|---|
| RRBS | Reduced representation bisulfite sequencing |
| WGBS | Whole-genome bisulfite sequencing |

| Tool | Reference/web-page |
|---|---|
| BSMap | https://code.google.com/archive/p/bsmap/ |
| BS-Seeker3 | https://github.com/khuang28jhu/bs3 |
| Bismark | https://github.com/FelixKrueger/Bismark |
| segemehl | http://www.bioinf.uni-leipzig.de/Software/segemehl/ |