

## 1 **Title**

2 **“Same Difference”**: Comprehensive evaluation of three DNA methylation measurement  
3 **platforms.**

## 4 **Authors**

5 Thadeous J Kacmarczyk<sup>1</sup>, Mame P. Fall<sup>1</sup>, Xihui Zhang<sup>1</sup>, Yuan Xin<sup>1</sup>, Yushan Li<sup>1</sup>, Alicia Alonso<sup>1</sup>,  
6 Doron Betel<sup>1,2</sup>

7 1. Department of Medicine, Division of Hematology/Oncology, Epigenomics Core Facility, Weill  
8 Cornell Medicine, New York, New York

9 2. Institute for Computational Biomedicine, Weill Cornell Medicine, New York, New York

10

11 Corresponding author: Thadeous J Kacmarczyk, [thk2008@med.cornell.edu](mailto:thk2008@med.cornell.edu)

## 12 **Abstract**

13 **Background:** DNA methylation in CpG context is fundamental to the epigenetic regulation of  
14 gene expression in high eukaryotes. Disorganization of methylation status is implicated in many  
15 diseases, cellular differentiation, imprinting, and other biological processes. Techniques that enrich  
16 for biologically relevant genomic regions with high CpG content are desired, since, depending on  
17 the size of an organism's methylome, the depth of sequencing required to cover all CpGs can be  
18 prohibitively expensive. Currently, restriction enzyme based reduced representation bisulfite  
19 sequencing and its modified protocols are widely used to study methylation differences. Recently,  
20 Agilent Technologies and Roche NimbleGen have ventured to both reduce sequencing costs and  
21 capture CpGs of known biological relevance by marketing in-solution custom-capture hybridization  
22 platforms. We aimed to evaluate the similarities and differences of these three methods

23 considering each targets approximately 10-13% of the human methylome.

24 **Results:** Overall, the regions covered per platform were as expected: targeted capture based  
25 methods covered >95% of their designed regions whereas the restriction enzyme-based method  
26 covered >70% of the expected fragments. While the total number of CpG loci shared by all  
27 methods was low, ~30% of any platform, the methylation levels of CpGs common across platforms  
28 were concordant. Annotation of CpG loci with genomic features revealed roughly the same  
29 proportions of feature annotations across the three platforms. Targeted capture methods  
30 encompass similar amounts of annotations with the restriction enzyme based method covering  
31 fewer promoters (~9%) and shores (~8%) and more unannotated loci (7-14%).

32 **Conclusions:** Although all methods are largely consistent in terms of covered CpG loci and cover  
33 similar proportions of annotated CpG loci, the restriction based enrichment results in more  
34 unannotated regions and the commercially available capture methods result in less off-target  
35 regions. Quality of DNA is very important for restriction based enrichment and starting material  
36 can be low. Conversely, quality of the starting material is less important for capture methods, and  
37 at least twice the amount of starting material is required. Pricing is marginally less for restriction  
38 based enrichment, and number of samples to be prepared is not restricted to the number of  
39 samples a kit supports. The one advantage of capture libraries is the ability to custom design  
40 areas of interest. The choice of the technique should be decided by the number of samples, the  
41 quality and quantity of DNA available and the biological areas of interest since comparable data  
42 are obtained from all platforms.

43

44 **Keywords:** Bisulfite sequencing, DNA methylation, Methylome capture, RRBS, 5mC, CpG

## 45 **Background**

46 DNA cytosine methylation in the form of 5-methylcytosine (5mC) in CpG context is an epigenetic  
47 marker that is important for regulation of gene expression. Changes in CpG methylation are  
48 implicated in many diseases, and proper methylation patterns are required for normal  
49 development [1]. Large scale studies such as ENCODE [2] and the Human Epigenomics  
50 Roadmap [3] have performed extensive profiling of 5mC in various cell lines and tissues revealing  
51 a rich and dynamic landscape of 5mC patterns in the human genome. Given the importance of  
52 these markers to cellular development and contribution to disease, a number of approaches have  
53 been developed for detecting the methylation status of cytosines [4], with bisulfite sequencing  
54 (BS-seq) being widely used to provide single-base quantitative measurement of 5mC (and 5-  
55 hydroxymethylcytosine, 5hmC). Bisulfite sequencing refers to massively parallel sequencing after  
56 chemical deamination of cytosines (C) to uracils (U), followed by polymerase chain reaction  
57 (PCR). The deamination of cytosines is accomplished by the use of sodium bisulfite, and this pre-  
58 treatment preserves both the methyl modification in 5mC and the 5-hydroxymethyl modification in  
59 5hmC [5]. The benchmark in methylome coverage is whole genome bisulfite sequencing (WGBS),  
60 which at 30X sequencing coverage, accounts for ~94% of all cytosines in the genome with 99.8%  
61 of them being CpG loci [6]. However, different WGBS library preparation protocols can bias region  
62 coverage. Since no method completely covers the methylome, and biologically relevant CpGs  
63 have been identified in known genomic features [1,7], developing focused assays considering  
64 these features is in demand with the caveat that these approaches will leave gaps in the  
65 methylome potentially excluding important CpGs.

66

67 There are several methods for acquiring DNA methylation levels and we investigated the  
68 characteristics of three currently widely used platforms: i) enrichment by enzymatic digestion

69 (MspI) enhanced reduced representation bisulfite sequencing (ERRBS)[8], ii) capture based  
70 Agilent SureSelect Methyl-Seq (SSMethylSeq), and iii) capture based Roche NimbleGen SeqCap  
71 Epi CpGiant (CpGiant).

72 In this paper we present an analysis of the methylation pattern of the human lung fibroblast IMR-  
73 90 cell line profiled by each of the platform protocols, using two technical replicate libraries for  
74 ERRBS, and two libraries each for SSMethylSeq and CpGiant, one at the manufacturer's  
75 suggested concentration and one at a reduced concentration (Table 1 and Additional file 1: Table  
76 S1). The capture libraries differ only in concentration of input material and are treated effectively  
77 as technical replicates. All libraries were sequenced to equivalent depth and compared to a library  
78 made using WGBS.

## 79 **Materials and Methods**

### 80 **Cell growth and DNA preparation**

81 IMR90 cells (American Type Culture Collection, Manassas, VA cat # CCL-186) were provided by  
82 Dan Hasson (Mount Sinai School of Medicine, New York, NY). DNA from  $5 \times 10^7$  cells was purified  
83 using the Gentra Puregene DNA kit according to manufacturer protocol (cat # 158389, Qiagen  
84 Valencia, CA). DNA was resuspended in TE, quantified using fluorometric quantification (Qubit 2.0  
85 ThermoFisher Scientific Waltham, MA) and quality was assessed by running on a 1% agarose gel.  
86

### 87 **ERRBS (digestion-based enhanced reduced representation bisulfite** 88 **sequencing)**

89 Two ERRBS libraries were prepared as described in Garrett-Bakelman, et al [8]. Briefly, 75 ng of  
90 DNA was digested with the methylation insensitive MspI enzyme (C<sup>^</sup>CGG). After end-repair, A-

91 tailing, and adapter ligation with Illumina TruSeq adapters, the region corresponding to 84-334bp  
92 was size-selected as two fractions. Each fraction was subjected to overnight bisulfite conversion  
93 (55 cycles of 95°C for 30 sec, 50°C for 15 min) using EZ DNA methylation kit (Cat # D5002, Zymo  
94 Research, Irvine CA). Purified bisulfite converted DNA was PCR amplified using TruSeq primers  
95 (Illumina Inc. San Diego, CA) for 18 cycles of denaturing, annealing and extension/elongation  
96 steps using Roche FastStart (cat # 03 553 361 001) – 94°C for 20 secs, 65°C for 30 secs, 72°C  
97 for 1 min, followed by 72°C for 3 min. The resulting libraries were normalized to 2nM and pooled  
98 at the same molar ratio. Samples were clustered at 6.5pM on a V3 paired-end read flow cell and  
99 sequenced for 100 cycles on an Illumina HiSeq 2500.

#### 100 **Agilent SureSelect Methyl-Seq (SSMethylSeq)**

101 Two libraries were made according to the company's specifications using 3ug and 1ug of DNA.  
102 DNA was sonicated using a Covaris S220 sonicator (Covaris, Woburn, MA) to obtain products of  
103 150-200bp. DNA was then end-repaired, A-tailed and ligated with methylated adapters to create  
104 pre-capture DNA libraries. DNA Libraries were then hybridized to the RNA SureSelect Human  
105 methyl-seq capture library at 65°C for 16 hrs. Hybridized products were purified by capture with  
106 Streptavidin beads and then subjected to bisulfite conversion (64°C for 2.5hr) using the Zymo EZ  
107 DNA Gold kit (Cat # D5005, Zymo Research, Irvine CA). The bisulfite treated libraries were PCR  
108 amplified for 8 cycles with Agilent Taq to get the required amount of DNA library and then indexed  
109 by another 6 cycles of PCR amplification to create the final libraries. Libraries were clustered at  
110 11pM on a V3 paired-end read flow cell and sequenced for 100 cycles on an Illumina HiSeq 2500.

#### 111 **Roche NimbleGen SeqCap Epi CpGiant (CpGiant)**

112 Two libraries were made using 1ug and 0.25ug of starting material, according to the company's  
113 specifications. DNA was sonicated using a Covaris S220 sonicator (Covaris, Woburn, MA) to

114 obtain products of 180-220bp. DNA was then end-repaired, A-tailed and ligated with methylated  
115 indexed-adapters to create pre-capture DNA libraries. The pre-capture libraries were bisulfite  
116 converted at 54°C for 1 hour using the Zymo EZ DNA Lightning kit (Cat # D5030, Zymo Research,  
117 Irvine CA). The bisulfite treated pre-capture libraries were PCR amplified with HiFi HotSart Uracil+  
118 polymerase (Cat# KK2802, Kapa Biosystems, Wilmington, MA). The amplified, bisulfite converted  
119 sample libraries were then hybridized to the probe pool of fully-, partially- and un-methylated  
120 cytosines from both strands of DNA oligos at 42°C for 72 hrs. Hybridized products were purified by  
121 capture with Capture Beads and PCR amplified for 15 cycles to create the final libraries. Libraries  
122 were clustered at 12pM on a V3 paired-end read flow cell and sequenced for 100 cycles on an  
123 Illumina HiSeq 2500.

## 124 **Whole Genome Bisulfite Sequencing (WGBS)**

125 Briefly, 100 ng of DNA were bisulfite converted using EZ DNA Methylation-Gold Kit (cat # D5005,  
126 Zymo Research Corporation, Irvine, CA) and the single stranded DNA obtained processed for  
127 library construction using the EpiGnome Methyl-Seq kit (Cat. # EGMK81324) as per manufacturer  
128 protocol (Illumina Madison, Madison, WI). The DNA was made double stranded by the use of  
129 5'tagged random hexamers and subsequently 3' tagged with terminal tagging oligo. The di-tagged  
130 DNA was enriched using 10 cycles of PCR, with PCR indexed-primers compatible with Illumina  
131 sequencing. The library was clustered at 7 pM on a V3 paired-end read flow cell and sequenced  
132 for 100 cycles on an Illumina HiSeq 2500.

## 133 **Computational analysis**

134 Illumina's CASAVA 1.8.2 was used to generate fastq files from basecalls. Raw fastq reads were  
135 processed by a custom pipeline that consists of: 1) filtering raw fastq reads for pass filter reads,  
136 2) trimming adapter sequence by FLEXBAR [9], 3) genomic alignments performed using Bismark

137 [10] and Bowtie [11] to reference human genome hg19, and 4) methylation calling by a custom  
138 PERL script [8]. Custom analysis scripts were written in R (version 3.3.0 [12]), including packages:  
139 Bioconductor - Biobase version 2.32.0 [13], GenomicRanges package version 1.24.0 [14],  
140 Beeswarm package version 0.2.3 [15], Affy package version 1.50.0 [16],  
141 BSgenome.Hsapiens.UCSC.hg19 package version 1.4.0 [17], PreprocessCore package version  
142 1.34.0 [18], UpSetR package version 1.2.0 [19], and VennDiagram package version 1.6.17 [20],  
143 were used to perform the analysis and are available at GitHub  
144 (<https://github.com/thk2008/methylseqplatformcomparison>). Sequencing data and methylation  
145 results are available to download from GEO (GSE83595)

## 146 **Results**

### 147 **Region coverage**

148 Targeted capture techniques (SSMethylSeq and CpGiant) have a designed set of genomic regions  
149 and therefore, a predicted set of CpGs covered. SSMethylSeq is specifically designed to capture  
150 CpGs from a single DNA strand, where the other platform capture CpGs from both DNA strands.  
151 The ERRBS protocol is considered targeted to the extent that the restriction digest produces  
152 consistent genomic fragments, sizes from 84-334bp are isolated during the library preparation gel  
153 extraction step. Note that since the DNA for WGBS is randomly sheared and coverage depends  
154 upon sequencing depth, there are no predicted regions per se and consequently WGBS is  
155 excluded from the region coverage analysis. Regions are considered overlapping if any two  
156 regions overlap by at least 10 bases. Overall, SSMethylSeq and MspI 84-334b have similar region  
157 length proportions, while CpGiant has fewer regions and more variable lengths (Figure 1A and  
158 Additional file 1: Table S2). Additionally, the two capture platforms cover similar regions, and  
159 similar amount of different regions from MspI 84-334b (Figure 1B).

160

## 161 **Description of analysis**

162 Even though the platforms are different, the data are all bisulfite sequencing data, thus the  
163 procedure for processing the data is the same. First we evaluated sequencing and alignment  
164 consistency. Next we compared strand symmetry of methylation values, since the SSMethylSeq  
165 platform is designed to cover predominantly one strand. Then we evaluated several measures to  
166 quantify properties of each platform relating to the target regions that are covered, methylation  
167 levels of cytosines in CpG context and coverage of genomic regions. Target region analysis  
168 includes the number of CpGs covered, the fraction of target regions covered, the coverage of  
169 target region CpGs, the overlap of CpGs across all platforms, and the concordance of methylation  
170 levels of overlapping CpGs. Similarly, genomic region evaluation is comprised of the number and  
171 fraction of annotated CpGs by genomic feature (i.e. CpG islands or shores – 2kb flanking the  
172 islands) and by gene part (i.e. promoter, exon, intron), and the overlap, coverage, and  
173 concordance of genomic annotations. We used the median absolute deviation (MAD), a robust  
174 measure of variability insensitive to outliers, to estimate the statistical dispersion in methylation  
175 level comparisons.

176

## 177 **Sample preparation and sequencing**

178 The general library preparation protocol for ERRBS is to digest input DNA with the methylation  
179 insensitive MspI enzyme followed by addition of adapters and barcodes, enrichment by gel size  
180 selection for fragments 84-334bp, bisulfite conversion, and amplification by PCR. The capture  
181 methods of SSMethylSeq and CpGiant are similar in that targeted capture of genomic regions is  
182 performed by hybridization to probes thereby providing a direct measure of CpG methylation at  
183 predefined regions. However, the methods differ in two important aspects. First, the SSMethylSeq



184 platform captures and measures methylation from only one DNA strand. Second, in CpGiant's  
 185 methodology bisulfite conversion is performed before hybridization to oligonucleotides whereas in  
 186 SSMethylSeq bisulfite conversion is performed prior to genomic capture. Hence, CpGiant's probe  
 187 design must capture the various combinations of DNA sequences that can result from bisulfite  
 188 conversion. In contrast to the targeted approaches WGBS is unbiased and relies on bisulfite  
 189 breakage of genomic DNA. All libraries were made from human lung fibroblast IMR90 cell line  
 190 DNA and prepared as described in the materials and methods with prominent differences outlined  
 191 in Table 1.

192 **Table 1: Protocols comparison**

	<b>ERRBS</b>	<b>SS-MethylSeq</b>	<b>CpGiant</b>	<b>WGBS</b>
DNA requirement	75ng >40kb	1, 3ug	0.25,1ug	100 ng >40kb
DNA processing	MspI digestion to completion followed by fractionation of 84-334bp	Sonication to 150-200bp	Sonication to 180-220bp	Single stranded and fragmented during bisulfite conversion
Enrichment method	Size fractionation of 84-334bp sizes	Hybridization to RNA capture library stranded design	Hybridization to oligo probes containing fully-, partially- and unmethylated cytosines from both strands	none
Bisulfite conversion step	Post-adaptor ligation	Post-hybridization capture	Pre-hybridization capture	Pre-adaptor ligation
Zymo Research Bisulfite Conversion kit	EZ-DNA Methylation (50°C, 55cycles)	EZ-DNA Methylation Gold (64°C, 2.5hr)	EZ-Methylation Lightning (54°C, 1hr)	EZ DNA Methylation-Gold
Total PCR amplification cycles	18 Post enrichment and bisulfite conversion	16 Post enrichment and bisulfite conversion (8 for amplification and 6 for Indexing)	28 12 post bisulfite conversion & 16 post enrichment	10 Post bisulfite conversion
DNA Polymerase (uracil tolerant)	FastStart Taq (Roche)	Taq2000 (Agilent Technologies)	HiFi HotSart Uracil+ (Kapa Biosystems)	FailSafe Enzyme (Epicentre)
Predicted number of targeted CpG	6.6M	3.7M	5.6M	56M

sites				
Relative price (library preparation and fixed coverage)	13.5 %	15.5%	15.5%	100%

193

## 194 **Sequencing and alignment characteristics**

195 Generally, all platforms produced similar sequencing results with no noticeable bias or reduced  
196 quality scores. The number of clusters and number of pass filter reads produced a typically  
197 consistent number of usable reads for all samples (254M +/- 37M), see Additional file 1: Table S3  
198 for more details.

199 Consistent with other bisulfite converted samples the number of uniquely mapped reads was ~  
200 70.8% sd=5.2% (Figure 2A). Across all platforms, a mean of 24.2%, s= 2.2% of reads were not  
201 aligned and ambiguously mapped reads showed a larger proportion for ERRBS (~12.1%), than  
202 WGBS (4.6%), SSMethylSeq (1.4%), or CpGiant (1.4%) (Figure 2A). These observations indicate  
203 the propensity of designed platforms to avoid repetitive regions or paralogous regions in contrast  
204 to ERRBS where no such selection is possible.

205

## 206 **Strand methylation symmetry**

207 Maintenance of symmetric methylation patterns across complementary CpG sites is required in  
208 order to preserve methylation patterns across cellular divisions. In principle, measuring  
209 methylation levels on one strand is sufficient to infer the cellular methylation state. Indeed, the  
210 SSMethylSeq platform is designed to capture only one strand in contrast to all other platforms that  
211 capture from both strands (Figure 2B). We note that this is different from measuring hemi-  
212 methylation states, which refers to different methylation between two parental alleles. This

213 requires genomic phasing information and is not considered in this analysis. To validate that  
214 methylation is indeed strand symmetric we compared the methylation values between  
215 complementary CpG sites (Figure 3). For ERRBS and CpGiant about 40% of the CpGs were  
216 complementary and analysis of SSMethylSeq was done with the ~2% of the probes that were  
217 complementary. We found strong agreement in methylation values in all samples and all platforms  
218 supporting the notion of symmetric methylation (Figure 3, mean MAD=0.28, sd=0.06). ERRBS  
219 contains a small set of discordant CpG sites (Figure 3A, B), where methylation values are  
220 inconsistent among complementary sites ( $\Delta > 99$ , 3.8% of sites ERRBS\_A, 1.0% in ERRBS\_B at  
221 10X coverage). This discordance is a consequence of the ERRBS library preparation where the  
222 MspI staggered cut sites are *de novo (in vitro)* filled with C' and G' to generate blunt ends. As a  
223 result, methylated CpGs at the restriction sites (i.e. at the ends of sequences fragments) are  
224 discordant. To correct this artificial discordance, it's common for single end sequencing to discard  
225 CpG values at MspI cut sites. We observed additional disparity occurring in paired-end  
226 sequencing where the sequenced fragments are shorter than the overall sequencing length where  
227 sequencing reads from opposite strands are overlapping. Since the fraction of discordant sites is  
228 small, we have kept them in the analysis. However, In order to maintain equal evaluation of CpG  
229 capture and methylation levels among all platforms, all subsequent analyses are based on CpG-  
230 units, which are defined as CpG sites with  $\geq 10$  spanning reads regardless of which strand the  
231 reads are mapped (i.e. if two CpG sites are complementary they are combined into one unit and  
232 then filtered for  $\geq 10x$  coverage).

233

### 234 **CpG-unit coverage and target region coverage.**

235 We evaluated the number of CpG-units covered and the mean coverage of all CpG-units that  
236 were sequenced at 10X depth or more. ERRBS, SSMethylSeq and CpGiant platforms cover on

237 average 3.8M CpG-units,  $sd=0.2M$  at mean coverage of 114.27,  $sd=23.6$  (Figure 4A,B and  
238 Additional file 1: Figures S1 and S2). In contrast, WGBS sequenced at ~400M paired-end reads  
239 covered 14.4M CpG-units (~50% of the 28M total CpG-units in the genome) at a mean coverage  
240 of 19x, demonstrating that achieving comparable coverage by WGBS is significantly more costly  
241 than targeted platforms. Therefore, these results confirm that while WGBS profiles a large portion  
242 of CpG-units, the targeted platforms provide a cost effective way to interrogate a limited, yet  
243 potentially most informative, set of CpG-units at considerably reduced cost.

244 Next, we evaluated the extent of coverage in the targeted regions and the fraction of their CpGs  
245 identified. Here, a region is considered as having coverage if at least one CpG-unit is sufficiently  
246 covered within the region. We observed that ERRBS covers on average 79.8% of its predicted  
247 regions 84-334bp, whereas CpGiant and SSMethylSeq cover, on average, 98.7% and 95.3% of  
248 their expected capture regions, respectively (Figure 4A). While region coverage and the coverage  
249 of CpG-units within the targeted regions are high, we observed that roughly 20%-40% of the CpG-  
250 units covered are outside the targeted regions (Figure 4A). In the case of ERRBS this is expected  
251 since restriction digestion and gel isolation can be variable. In the case of capture methods, this  
252 may indicate either cross-hybridization of the probes to other genomic locations or hybridization of  
253 longer fragments. Since there are typically several CpG-units within a region, we looked at the  
254 distribution of the fraction of CpG-units covered per region and found that nearly all CpG-units in  
255 targeted regions are covered although for both SSMethylSeq and ERRBS a number of CpGs in  
256 targeted regions are not represented presumably due to coverage bias (Figure 4C).

## 257

### 258 **CpG-unit overlap and methylation levels concordance**

259 Overall, the average number of common CpG-units covered by all three platforms is ~30% +/- 1%  
260 (Figure 4D). When comparing the pairwise overlap of shared CpG-units across all samples, we

261 observe high overlap between intra-platform replicates with median number of shared CpG-units  
262 ~3.58M, median overlap 92.1% and high methylation level concordance with mean MAD = 0.23  
263 and sd = 0.05 (Additional file 1: Figure S3 and Table S4). Among inter-platform samples we  
264 observed less overlap; median number of shared CpG-units ~1.5M, median overlap 39.1%  
265 (Figure 5 upper triangle, Additional file 1: Figure S3 and Table S4). Methylation levels of common  
266 CpG-unit's across all platforms are highly concordant indicating that methylation level  
267 measurements are consistent and reproducible with average MAD = 0.41 and sd = 0.25 (Figure 5  
268 lower triangle, Additional file 1: Figure S3 and Table S4). Inter-platform concordance was slightly  
269 lower than intra-platform concordance with mean MAD = 0.46 and sd = 0.27 (Figure 5 lower  
270 triangle, Additional file 1: Figure S3 and Table S4). These results demonstrate that while the  
271 platforms differ in their capture approaches and bisulfite conversion kits (Table 1) these  
272 differences are not biasing methylation measurements. The differences among the platforms,  
273 therefore, are largely restricted to variations in targeted regions and not in methylation  
274 measurements.

275

## 276 **Coverage of genomic feature regions by CpG-units**

277 Next we were interested in what genomic features are covered by each platform and the degree of  
278 coverage. Analogous to the previous analysis, here we define a region as being one of exon,  
279 intron, promoter, CpGi, or shore, and coverage is a genomic region that contains at least one  
280 detected CpG-unit. It should be noted that the genomic feature annotations are not mutually  
281 exclusive and that some CpG-units are annotated by more than one category. Naturally, the  
282 targeted platforms focus on genomic regions known to play important roles in epigenetic  
283 regulation such as promoter regions, CpG islands, and CpG shores and ERRBS covers the same  
284 regions, but to a lesser degree (Figure 6A). Moreover, each platform covers similar proportions of

285 CpG-units in particular genomic region (Figure 6B). Conversely, we looked at the number of CpG-  
286 units that have an annotation and observed similar trends for the three platforms with ERRBS  
287 having a larger proportion of unannotated CpG-units (~27%, Figure 6C). Overall, no platform is  
288 significantly enriched for particular genomic region although ERRBS is slightly less represented in  
289 promoters, CpG islands and CpG shores, while having more representation of unannotated CpG-  
290 units.

291

## 292 **Overlap of CpG-units annotated with a genomic feature**

293 We evaluated the overlap of annotated CpG-units to determine whether any platform is enriched  
294 for a particular genomic feature (e.g. exons, introns, promoters, CpG islands and shores). Again, it  
295 should be mentioned that the genomic feature annotations are not mutually exclusive and that  
296 some CpG-units are annotated by more than one category. However, a CpG-unit is counted as  
297 annotated if it has one or more annotations. Comparing the overlap in annotations across  
298 platforms, we see a similar grouping as above with the methylation levels, intra-platform overlap is  
299 high (mean overlap 93.7% sd=2.9%), and inter-platform overlap is lower (mean 55.1% sd=15.5%)  
300 (Figure 7 and Additional file 1: Figures S4-S9). Thus, we observed similar proportions of genomic  
301 region coverage across all platforms, but lower overlap of annotated CpG-units suggesting  
302 coverage of different loci within those regions.

303

## 304 **Discussion**

305 We performed a systematic analysis of the characteristics of three prominent platforms for  
306 measuring DNA methylation levels: ERRBS, Agilent SureSelect Methyl-Seq (SSMethylSeq),  
307 Roche NimbleGen SeqCap Epi CpGiant (CpGiant), and WGBS with the goal of identifying

308 whether one method outperforms the others for any particular property.

309 We assessed the expected region coverage for SSMethylSeq and CpGiant and MspI (for ERRBS)

310 and observed that, overall, the CpG sites covered by each platform are largely distinct from each

311 other. SSMethylSeq and CpGiant cover roughly 70% of the same CpGs, while MspI has only

312 ~30% overlap with either capture. However, each platform covers a large fraction of its targeted

313 regions. SSMethylSeq and CpGiant cover greater than 95% of their designed capture regions and

314 ERRBS covers at least 77% of its expected fragments. Furthermore, methylation levels of

315 overlapping CpG-units are highly concordant.

316 Finally, each platform covered roughly the same proportions of genomic features (CpG islands,

317 shores, promoters, exons, and introns), but profiled different CpG sites within those regions. In

318 addition to the differences in the targeted CpG positions, there are also differences in the library

319 preparation protocols. ERRBS can be performed using small amounts of starting material (as little

320 as 75ngs), whereas SSMethylSeq requires 1ug and CpGiant 0.25ug of starting DNA (although

321 those are expected to reduce with further optimization by the vendors). ERRBS, or digestion-

322 based methods, can cover certain genomic regions that are not amenable to capture. ERRBS

323 provides comparable data to the other platforms, although there is more variability among the

324 profiled CpGs. Capture platforms are more precise and can be customized for profiling specific

325 genomic regions of interest. In addition, capture platforms are the only methods available for

326 methylation profiling of low quality, degraded DNA although library preparation is more labor

327 intensive and efficiency of capture depends on shearing.

## 328 **Conclusions**

329 Epigenetic state is a fundamental element of cellular development and regulation. Therefore,

330 accurate, reproducible and cost effective approaches for profiling DNA methylation are important

331 for advancements in biomedical research. While the cost of sequencing continues to decrease,

332 reaching sufficient coverage for reliable measurement of CpG methylation by WGBS is still  
333 prohibitive. We conclude from our comparative study that capture-based approaches provide  
334 comparable results and cover more precisely their intended targets than ERRBS, which is a  
335 restriction enzyme based approach. They also provide the added flexibility of designing custom  
336 capture for surveying regions of interest. On the other hand, digestion based protocols are  
337 currently more cost effective and may be the only option for clinical samples where the amount of  
338 input DNA is limited. In the absence of any prior knowledge about the genomic regions of interest  
339 for a particular study, the choice of methylation profiling platform should be guided by cost and  
340 amount and quality of input DNA.

341

## 342 **List of abbreviations**

343 5mC: 5-methylcytosine, BS-seq: bisulfite sequencing, CpGiant: Roche NimbleGen SeqCap Epi  
344 CpGiant, CpGi: CpG island, ERRBS: enhanced reduced representation bisulfite sequencing,  
345 GEO: gene expression omnibus, MAD: median absolute deviation, PCR: polymerase chain  
346 reaction, SSMethylSeq: Agilent SureSelect Methyl-Seq, WGBS: whole-genome bisulfite  
347 sequencing

348

## 349 **Declarations**

350 *Ethics approval and consent to participate*

351 Not applicable

352 *Consent for publication*

353 Not applicable



354 *Availability of data and material*

355 The data that support the findings of this study are available in NCBI GEO  
356 (<http://www.ncbi.nlm.nih.gov/geo/>) under Accession Number GSE83595. The code for the analysis  
357 is available from the code hosting platform GitHub  
358 (<https://github.com/thk2008/methylseqplatformcomparison>)

359 *Competing interests*

360 The authors declare that they have no competing interests.

361 *Funding*

362 DB is funded by grants from the Starr Consortium (I8-A8-132) and Tri-SCI (2013-036).

363 *Authors' contributions*

364 DB, AA Conceived and designed the experiments. MPF, XZ, YL performed the experiments. TJK  
365 performed the bioinformatic analysis and prepared the figures. TJK, DB, AA, analyzed the data.  
366 TJK, DB drafted the manuscript. TJK, DB, AA, MPF reviewed the manuscript. All authors read and  
367 approved the final manuscript.

368 *Acknowledgements*

369

370

371 *Additional file*

372 Additional file1. **Table S1.** Library input details. **Table S2.** Target region properties and CpGs  
373 covered. **Table S3.** Sequencing details. **Figure S1.** Number of CpG-units covered. **Figure S2.**

374 Mean and median coverage per CpG-unit. **Figure S3**. Intra- and Inter- platform CpG-unit overlap  
375 and methylation levels concordance. **Table S4**. Intra- and Inter- platform details. **Figure S4**.  
376 Overlap of exon annotation of CpG-units as UpSet plot. **Figure S5**. Overlap of intron annotation of  
377 CpG-units as UpSet plot. **Figure S6**. Overlap of promoters annotation of CpG-units as UpSet plot.  
378 **Figure S7**. Overlap of CpG island annotation of CpG-units as UpSet plot. **Figure S8**. Overlap  
379 CpG shores annotation of CpG-units as UpSet plot. **Figure S9**. Overlap of unannotated CpG-units  
380 as UpSet plot.

381

## 382 **References**

- 383 1. Jones PA. Functions of DNA methylation: islands, start sites, gene bodies and beyond. *Nat*  
384 *Rev Genet.* 2012;13(7):484–92.
- 385 2. The ENCODE Project Consortium. The ENCODE (ENCyclopedia Of DNA Elements)  
386 Project. *Science.* 2004;306(5696):636–40.
- 387 3. Consortium RE, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, et al. Integrative  
388 analysis of 111 reference human epigenomes. *Nature.* 2015;518(7539):317–30.
- 389 4. Harris RA, Wang T, Coarfa C, Nagarajan RP, Hong C, Downey SL, et al. Comparison of  
390 sequencing-based methods to profile DNA methylation and identification of monoallelic  
391 epigenetic modifications. *Nat Biotechnol.* 2010;28(10):1097–105.
- 392 5. Frommer M, McDonald LE, Millar DS, Collis CM, Watt F, Grigg GW, et al. A genomic  
393 sequencing protocol that yields a positive display of 5-methylcytosine residues in individual  
394 DNA strands. *Proc Natl Acad Sci USA.* 1992;89(5):1827–31.
- 395 6. Lister R, Pelizzola M, Dowen RH, Hawkins RD, Hon G, Tonti-Filippini J, et al. Human DNA  
396 methylomes at base resolution show widespread epigenomic differences. *Nature.*  
397 2009;462(7271):315–22.
- 398 7. The BLUEPRINT consortium. Quantitative comparison of DNA methylation assays for  
399 biomarker development and clinical applications. *Nat Biotech.* 2016;34:726–737
- 400 8. Garrett-Bakelman FE, Sheridan CK, Kacmarczyk TJ, Ishii J, Betel D, Alonso A, et al.  
401 Enhanced reduced representation bisulfite sequencing for assessment of DNA methylation  
402 at base pair resolution. *J Vis Exp.* 2015;(96):e52246.
- 403 9. Dodt M, Roehr JT, Ahmed R, Dieterich C. FLEXBAR-Flexible Barcode and Adapter

- 404 Processing for Next-Generation Sequencing Platforms. *Biology (Basel)*. 2012;1(3):895–905.
- 405 10. Krueger F, Andrews SR. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq  
406 applications. *Bioinformatics* . 2011;27(11):1571–2.
- 407 11. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of  
408 short DNA sequences to the human genome. *Genome Biol* . 2009;10(3):R25.
- 409 12. R Core Team. R: A Language and Environment for Statistical Computing . Vienna, Austria;  
410 2016. Available from: <https://www.r-project.org/>
- 411 13. Huber, W., Carey, J. V, Gentleman, R., et al. Orchestrating high-throughput genomic  
412 analysis with Bioconductor. *Nat Methods* . 2015;12(2):115–21.
- 413 14. Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, et al. Software for  
414 Computing and Annotating Genomic Ranges. *{PLoS} Comput Biol* . 2013;9(8).
- 415 15. Eklund A. beeswarm: The Bee Swarm Plot, an Alternative to Stripchart . 2016. Available  
416 from: <https://cran.r-project.org/package=beeswarm>
- 417 16. Gautier L, Cope L, Bolstad BM, Irizarry RA. affy--analysis of Affymetrix GeneChip data at the  
418 probe level. *Bioinformatics*. 2004;20(3):307–15.
- 419 17. Team TBD. BSgenome.Hsapiens.UCSC.hg19: Full genome sequences for Homo sapiens  
420 (UCSC version hg19). 2014.
- 421 18. Bolstad BM. preprocessCore: A collection of pre-processing functions. 2016. Available from:  
422 <https://github.com/bmbolstad/preprocessCore>
- 423 19. Gehlenborg N. UpSetR: A More Scalable Alternative to Venn and Euler Diagrams for  
424 Visualizing Intersecting Sets. 2016. Available from: [https://cran.r-](https://cran.r-project.org/package=UpSetR)  
425 [project.org/package=UpSetR](https://cran.r-project.org/package=UpSetR)
- 426 20. Chen H. VennDiagram: Generate High-Resolution Venn and Euler Plots. 2016. Available  
427 from: <https://cran.r-project.org/package=VennDiagram>
- 428

## 429 **Figure legends**

### 430 **Figure 1. Length and overlap of design regions and MspI predicted regions.**

431 **A)** Boxplot showing the distributions of targeted region's lengths. SSMethylSeq and MspI 84-334b  
432 (ERRBS) show similar region proportions, while CpGiant regions are generally longer and are  
433 more variable in length. **B)** Venn diagrams showing the pairwise overlap comparison of each of  
434 the platform's regions and the total number of regions. Circle size is proportional to the size of the

435 set. The two capture platforms have a higher degree of common regions with each other than  
436 either one with MspI 84-344b.

437

438 **Figure 2. Read alignment and strand parity.**

439 **A)** Percent alignment of uniquely aligned reads (green), ambiguously mapped reads (grays),  
440 reads with no alignment (pink) and rejected reads (blue). **B)** The fraction of CpG's covered at  
441  $\geq 10x$  coverage grouped by strand, plus strand (blue) and the minus strand (pink). The strand  
442 specific protocol of SSMethylSeq platform is evident by the high proportion of reads mapped to  
443 the negative strand.

444

445 **Figure 3: Strand symmetry of methylation values MA-plots.**

446 MA-plots of the log average of the methylation levels (A) on the x-axis and log ratio of the  
447 methylation levels (M) on the y-axis, between complementary CpG positions. Median absolute  
448 deviation (MAD) values are used to evaluate the agreement in methylation levels. The bimodal  
449 nature of methylation patterns (mostly unmethylated or methylated) is reflected in the high density  
450 at both ends of the x-axis. The artificially discordant sites introduced during ERRBS library  
451 preparation are identified as increased density off the center line at the low methylation values (at  
452  $A < 0$  range) in panels A, B. Panels are: **A)** ERRBS\_A, **B)** ERRBS\_B, **C)** SSMethylSeq\_A, **D)**  
453 SSMethylSeq\_B, **E)** CpGiant\_A, **F)** CpGiant\_B, and **G)** WGBS.

454

455 **Figure 4. Platform CpG-unit region coverage and CpG-unit overlap.** **A)** Number of CpG-units  
456 identified in targeted and off target regions by each platform. WGBS (orange) covers ~14.4M  
457 CpG-units however there is no notion of targeted regions in this platform. The targeted platforms  
458 predicted total CpG-units are depicted as gray bar and coverage of the CpG-units in the predicted  
459 regions (a.k.a “on-target”) are shown in blue. CpG-units outside the predicted set (off-target) are

460 shown in red bars. The red square points and right scale represent the percent recovery of  
461 targeted CpG-units. **B)** Distribution of coverage per CpG-unit. **C).** Density plots of the percent of  
462 the dataset's regions (color scale) with the fraction of CpG-units covered per region. These plots  
463 demonstrate that the reduced number of recovered CpG-units in ERRBS relative to the other  
464 platforms (red squares in panel A) are attributed to increased number of missed CpG-units shown  
465 as increased density at 0 region. **D)** Triple Venn diagram showing the proportion of overlap of  
466 CpG-units for the three platforms; size is not proportional to dataset.

467

468 **Figure 5. Inter-platform CpG-unit overlap and methylation levels concordance.** The upper  
469 triangle shows the total number of CpG-units (blue) and number of overlapping CpG-units  
470 between any two samples (red) with the percent overlap indicated. The lower triangle shows MA-  
471 plots of common CpG-units between any two samples, where M is the log ratio of the methylation  
472 levels and A is the average log methylation level from the two platforms. There are blue clouds  
473 that at log scale show the variance in methylation at low levels. See Additional file 1: Figure S3  
474 for all pairwise comparisons.

475

476 **Figure 6: Summary of coverage and representation of annotated genomic regions by each**  
477 **platform. A)** The fraction of genomic feature covered by sample CpG-units where at least one  
478 CpG-unit is covered in a region. **B)** The fraction of a region's CpG-units covered, and **C)** The  
479 fraction of CpG-units annotated with a genomic feature. CpGi (CpG island).

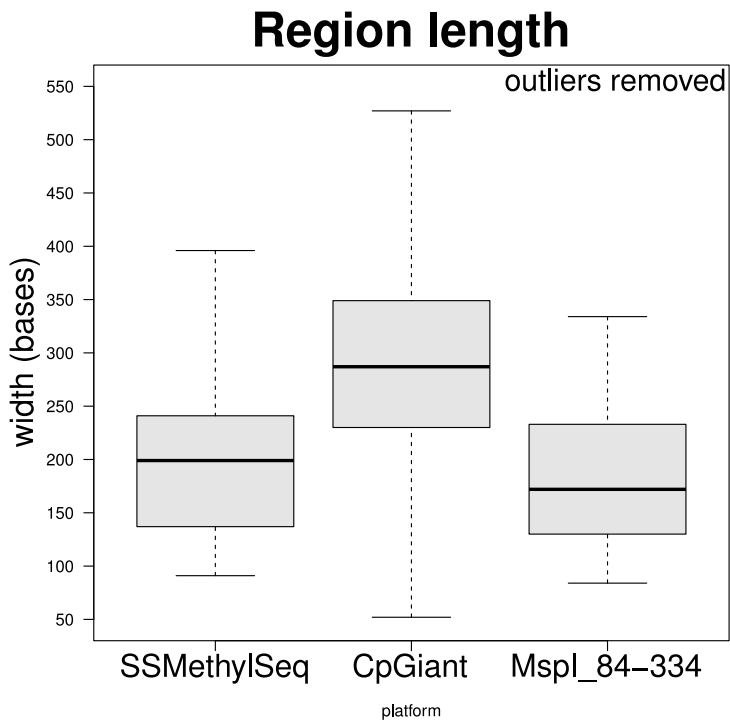
480

481 **Figure 7: CpG-unit annotations overlap.** Overlaps of all CpG-unit annotations across all  
482 platforms. The UpSet visualization technique [19] for set intersections is displayed as matrix  
483 layout. Horizontal bars on the lower left indicate the total number of annotated CpG-units in the  
484 set. Dark circles in the matrix (lower right) indicate sets that are part of the intersection. Bars in the

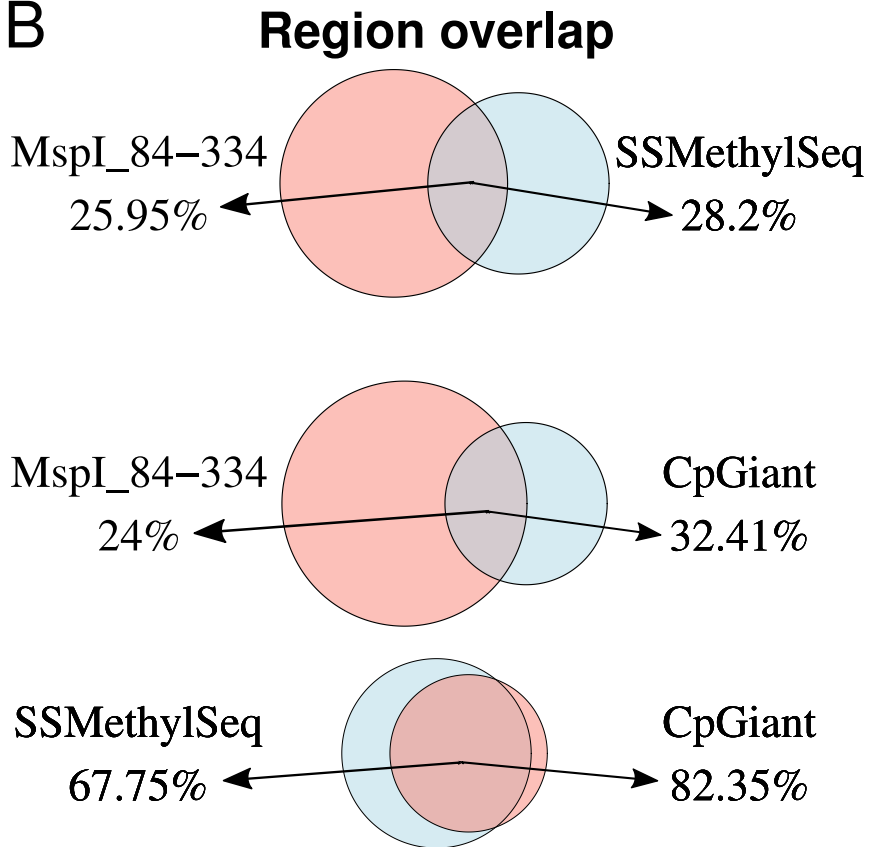
485 main plot area (upper right) indicate the number of intersecting CpG-units for the sets represented  
486 by the dark circles. Roughly 33% of the annotations are common to all three suggesting that while  
487 there may be similar proportions of CpG-unit annotations (i.e. they may be covering similar  
488 regions), they are covering different loci within those regions.

# Figure 1

## A



## B



# Figure 2

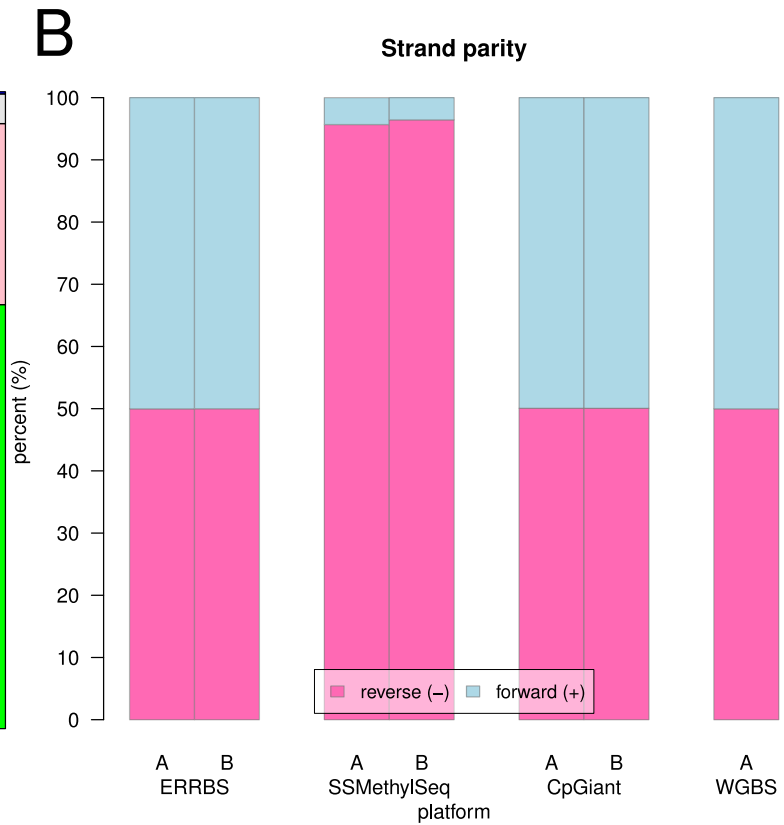
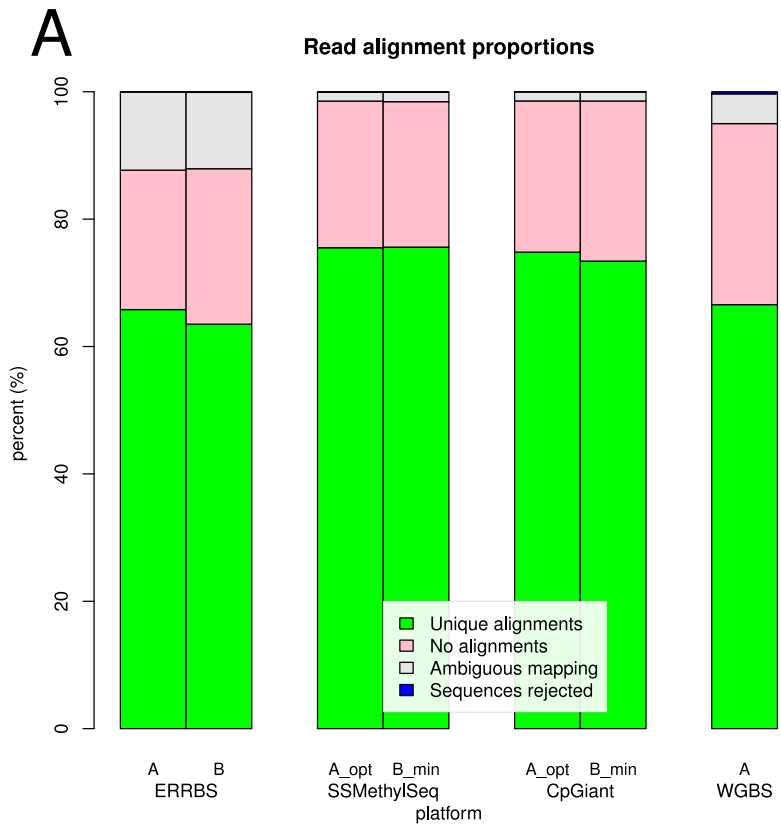
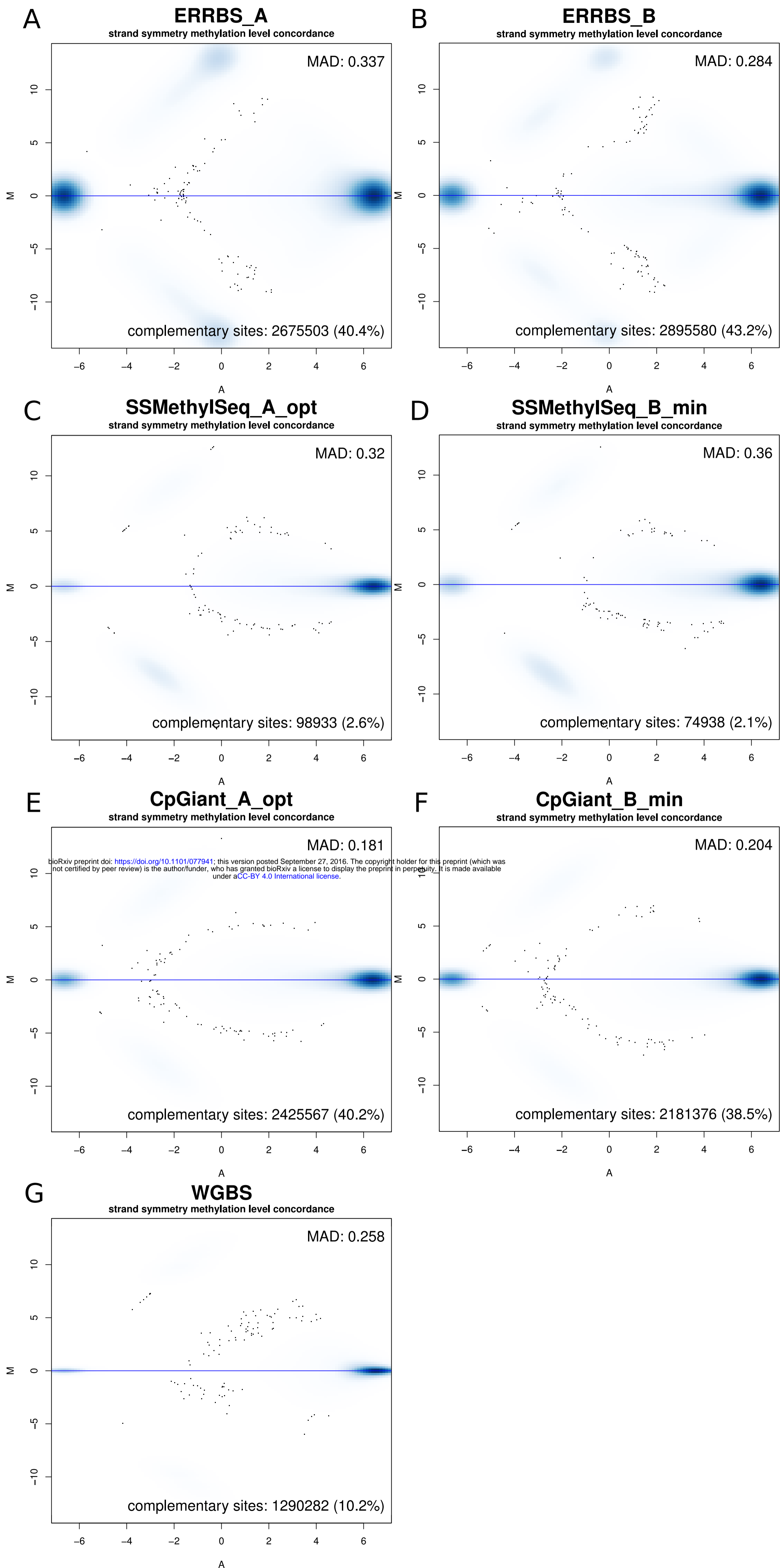




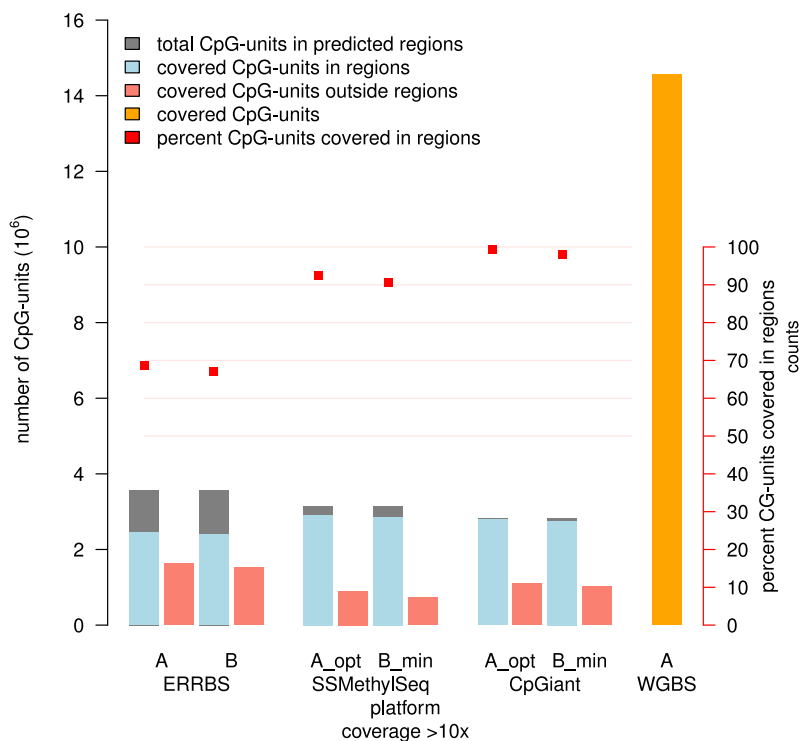
Figure 3



# Figure 4

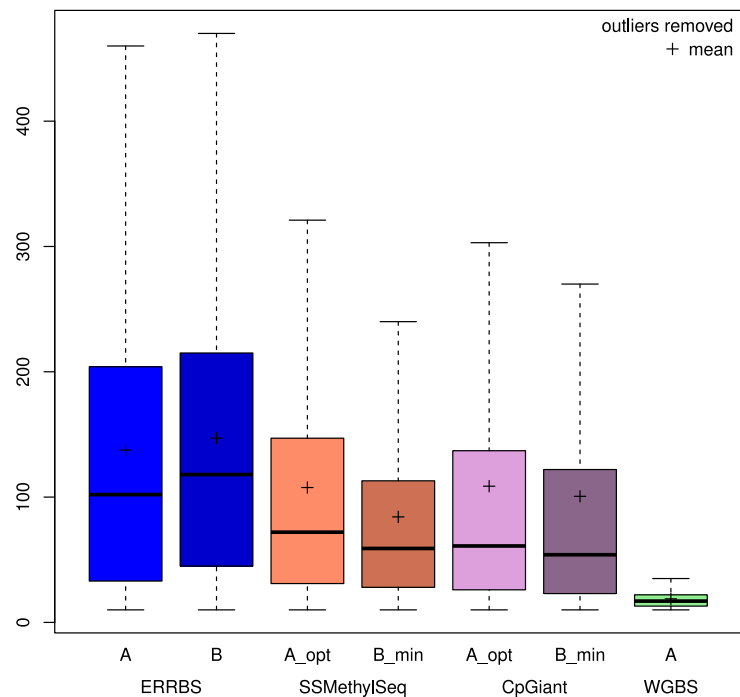
## A

### Number of CpG-units



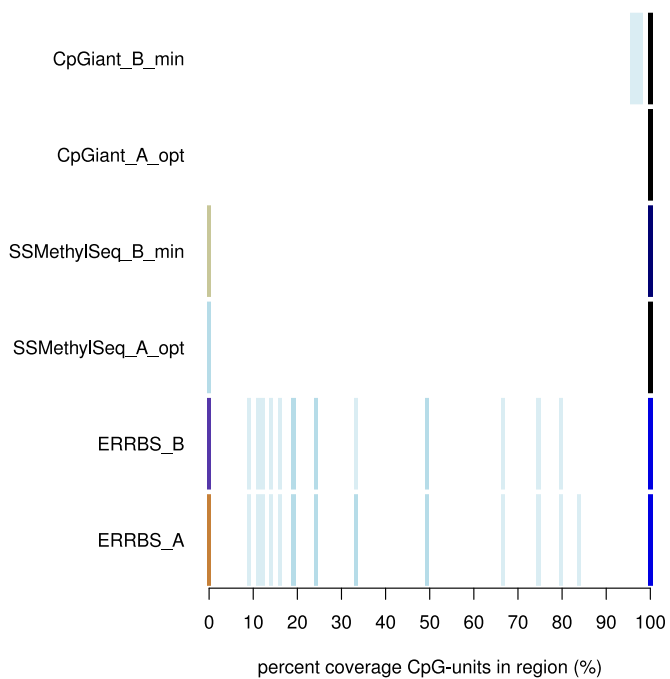
## B

### CpG-unit coverage $\geq 10x$



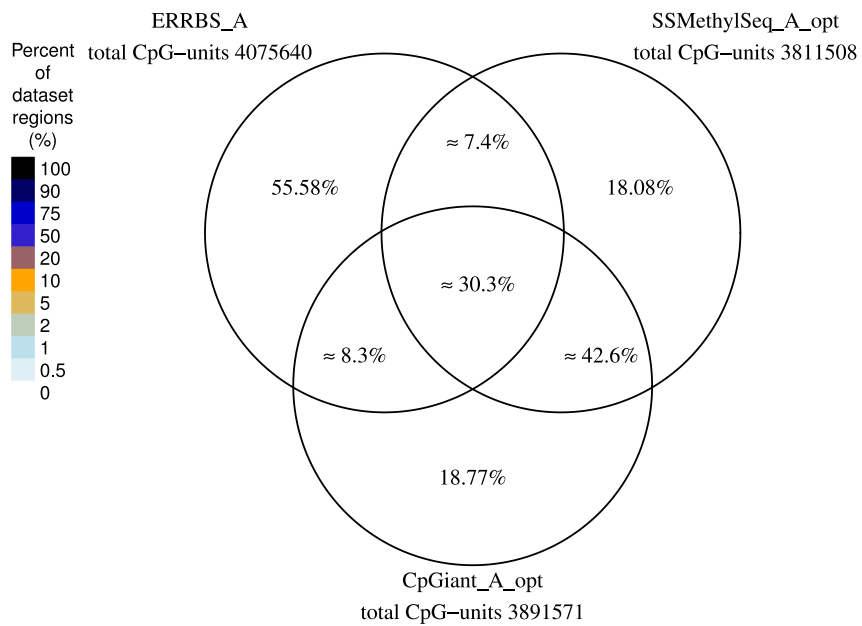
## C

### Percent of target regions with percent CpG-unit coverage



## D

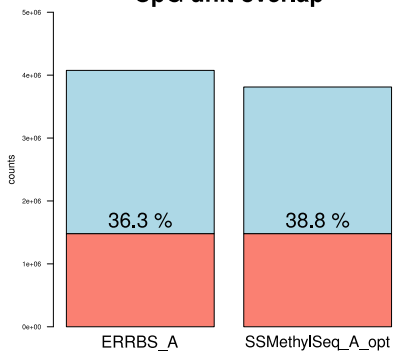
### CpG-unit overlap



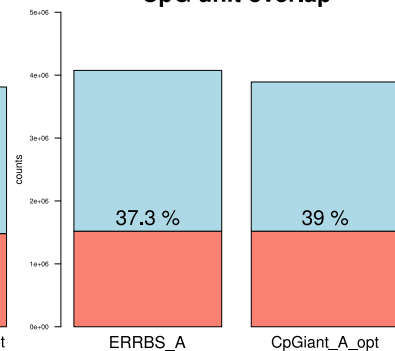
# Figure 5

ERRBS\_A

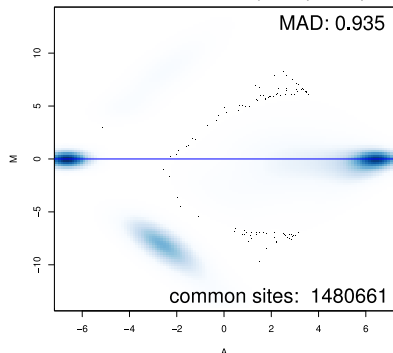
CpG unit overlap



CpG unit overlap

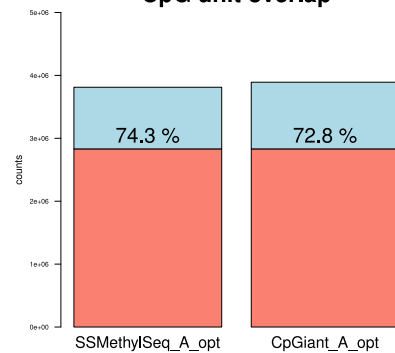


ERRBS\_A::SSMethylSeq\_A\_opt

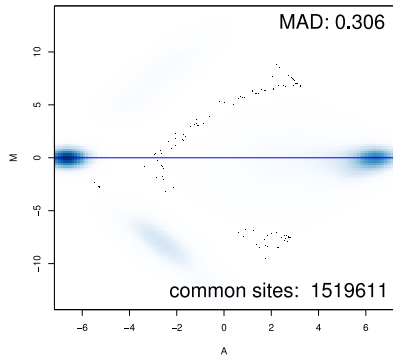


SSMethylSeq\_A

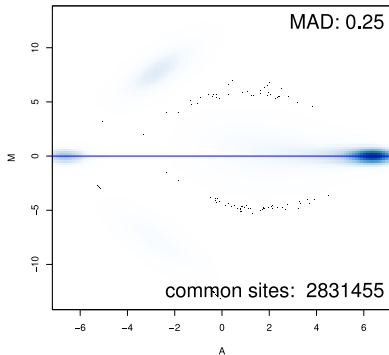
CpG unit overlap



ERRBS\_A::CpGiant\_A\_opt



SSMethylSeq\_A\_opt::CpGiant\_A\_opt

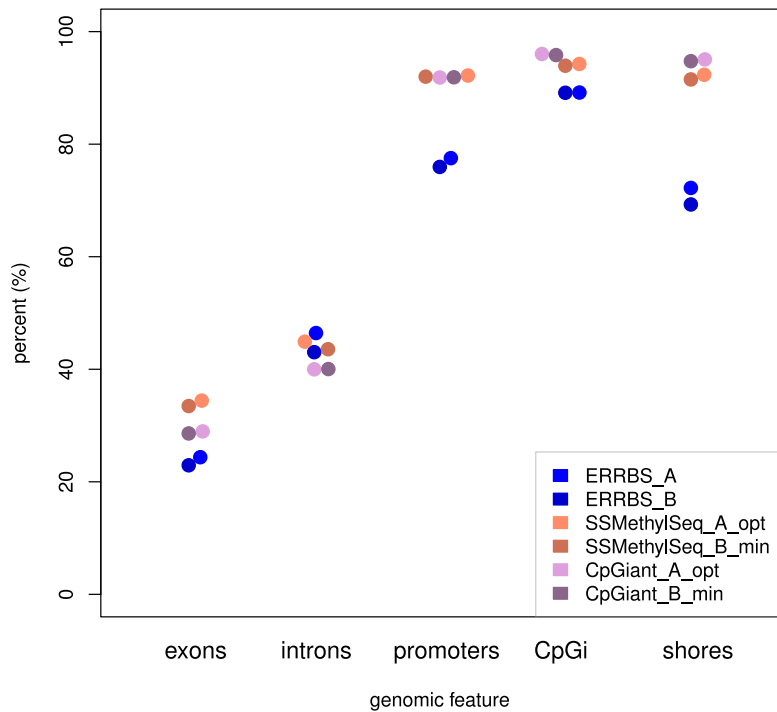


CpGiant\_A

# Figure 6

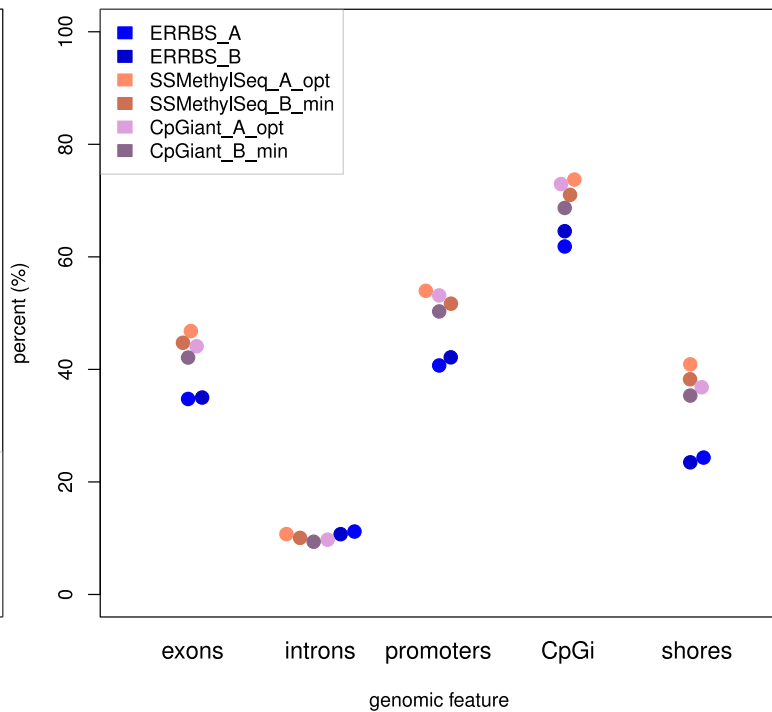
## A

### Region coverage



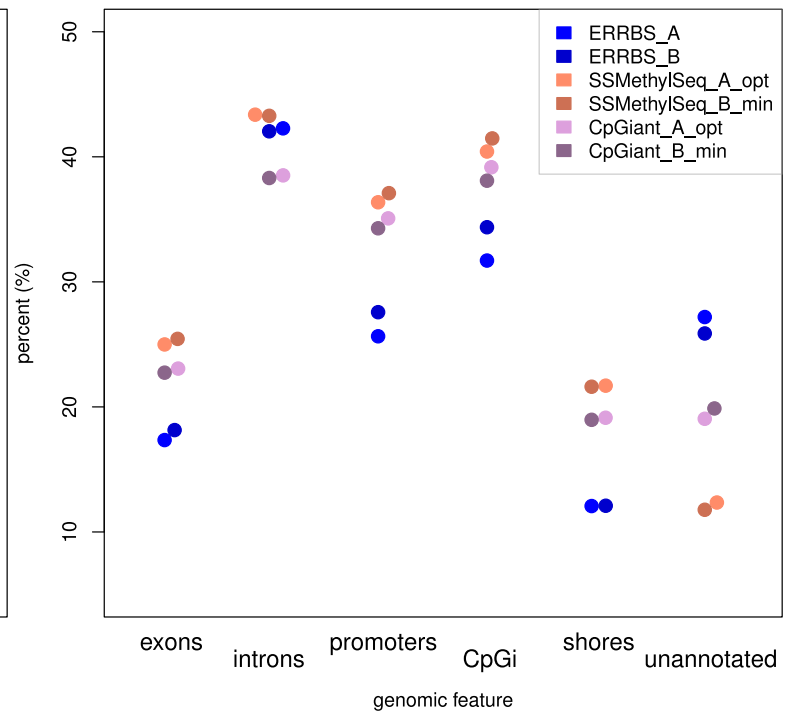
## B

### Region CpG-unit coverage



## C

### Sample CpG-units annotations



# Figure 7

Annotation overlap

