

MC profiling: a novel approach to analyze DNA methylation heterogeneity in genome-wide bisulfite sequencing data

Giulia De Riso^{1,*}, Antonella Sarnataro¹, Giovanni Scala², Mariella Cuomo^{1,3}, Rosa Della Monica³, Stefano Amente¹, Lorenzo Chiariotti^{1,3}, Gennaro Miele^{4,5} and Sergio Coccoza¹

¹Department of Molecular Medicine and Medical Biotechnology, University of Naples Federico II, via Sergio Pansini 5, 80131 Naples, Italy, ²Department of Biology, University of Naples Federico II, Via Vicinale Cupa Cintia 21, 80126 Naples, Italy, ³CEINGE - Biotechnologie Avanzate, Via Gaetano Salvatore, 486, 80145 Naples, Italy, ⁴Department of Physics “E. Pancini”, University of Naples “Federico II”, Via Cinthia, 80126 Naples, Italy and ⁵Istituto Nazionale di Fisica Nucleare (INFN), Sezione di Napoli, 80126 Naples, Italy

Received September 13, 2022; Revised November 24, 2022; Editorial Decision December 01, 2022; Accepted December 08, 2022

ABSTRACT

DNA methylation is an epigenetic mark implicated in crucial biological processes. Most of the knowledge about DNA methylation is based on bulk experiments, in which DNA methylation of genomic regions is reported as average methylation. However, average methylation does not inform on how methylated cytosines are distributed in each single DNA molecule. Here, we propose Methylation Class (MC) profiling as a genome-wide approach to the study of DNA methylation heterogeneity from bulk bisulfite sequencing experiments. The proposed approach is built on the concept of MCs, groups of DNA molecules sharing the same number of methylated cytosines. The relative abundances of MCs from sequencing reads incorporates the information on the average methylation, and directly informs on the methylation level of each molecule. By applying our approach to publicly available bisulfite-sequencing datasets, we individuated cell-to-cell differences as the prevalent contributor to methylation heterogeneity. Moreover, we individuated signatures of loci undergoing imprinting and X-inactivation, and highlighted differences between the two processes. When applying MC profiling to compare different conditions, we identified methylation changes occurring in regions with almost constant average methylation. Altogether, our results indicate that MC profiling can provide useful insights on the epigenetic status and its evolution at multiple genomic regions.

INTRODUCTION

DNA methylation is a heritable epigenetic mark consisting in the enzyme-mediated addition of a methyl-group to deoxyribonucleotides (1–3). In mammals, DNA methylation mainly involves cytosines in CpG context (1–3). DNA methylation has been shown to regulate gene expression and genome stability, and has been implicated in crucial biological processes, like genomic imprinting and X inactivation (3–6). In cell differentiation, DNA methylation shapes fate and engraves the identity of cells (1,7). Its dysregulation has been linked to plenty of pathological conditions (8–11).

Several experimental techniques have been developed to study DNA methylation (12). Among them, bisulfite sequencing techniques are widely adopted to assess the methylation status at single base resolution, either at targeted regions or at genome-wide level (12–16).

Single base DNA methylation is usually reported as the fraction of molecules in which a given cytosine is methylated (17). Genome-wide methylation analysis have highlighted that most cytosines are not evenly methylated in different molecules (18). Cellular heterogeneity and allele specific methylation are potential sources of this molecular heterogeneity (19).

Evidence has been provided that DNA methylation is regulated in larger genomic regions, with sets of neighboring cytosines working as functional units (20–23). DNA methylation analysis has indeed turned to the study of the average methylation of DNA regions (the fraction of methylated cytosines in a given region), and on the identification of regions with consistently different DNA methylation levels between groups of samples (differentially methylated regions, DMR) (17). Most of the current knowledge on DNA methylation and its implication in

*To whom correspondence should be addressed. Tel: +39 0817463757; Email: giulia.deriso@unina.it

health and disease status is founded on this latter approach (24–26).

However, the overall average methylation of a region does not inform on how this amount is contributed by the average methylation of single DNA molecules. As an example, an average methylation value of 0.5 for a given locus could result from a homogenous pool of half methylated molecules, or from an heterogeneous, balanced set composed of fully methylated and unmethylated molecules, or even from more heterogeneous pools (Figure 1 of (27)).

Single-cell DNA methylation assays have highlighted extensive cell-to-cell differences in regional DNA methylation (28,29), and have demonstrated that cellular heterogeneity can have a functional impact. For example, epigenetic variability at regulatory elements has been linked with gene expression variability (30,31). However, single-cell DNA methylation assays are still limitedly adopted due to the high cost and large sparsity of produced data (28,32).

Besides single cell techniques, analysis of DNA methylation patterns in bisulfite sequencing reads has also been adopted to analyze DNA methylation heterogeneity in bulk samples (19,28,33–36). In this context, mathematical modeling has been applied to estimate the distribution of methylation levels from Whole Genome Bisulfite Sequencing (WGBS) data, giving insights on its disposition across the genome, its evolution upon differentiation, aging and cancer, and its relationship with the genetic background (37–39).

In previous studies high-coverage amplicon bisulfite sequencing allowed us to directly estimate the distribution of methylation levels from supporting sequencing reads at targeted regions. Our approach, here referred as MC profiling, was based on the concept of Methylation Classes (MCs), i.e. groups of molecules holding the same amount of methylated cytosines, and allowed us to gain insights on the regulatory mechanisms of DNA methylation (40,41).

In this study, we extended MC profiling to genome-wide bisulfite sequencing data, with the aim to explore DNA methylation heterogeneity at multiple target regions.

In this setting, MC profiling identified cell-to-cell differences as the prevalent contributor to DNA methylation heterogeneity, with allele differences emerging in a small fraction of analyzed regions. Moreover, MC profiling led to the identification of signatures of loci undergoing genomic imprinting and X inactivation, and highlighted differences between the two processes. When applied to a dynamic system, MC profiling identified DNA methylation changes in regions with almost constant average methylation. Altogether, our results indicate that MC profiling can provide useful insights on the epigenetic status and its evolution at multiple genomic regions.

MATERIALS AND METHODS

MC profiling

Epilocus definition and MC profiling. We analyzed the methylation status of regions holding 4 CpG sites, with the first and the fourth CpGs delimiting the region. We refer to these regions as epiloci in the following text (Figure 1A). For a certain epilocus, we analysed the abundance of each possible MC class, i.e. groups of DNA molecules bearing

the same number of methylated cytosines. For a locus bearing four CpG sites, the number of MC classes is equal to 5. We depicted the DNA methylation status of a given epilocus through its inherent MC profile, i.e. the set of the relative abundances of the 5 MCs. To compute this MC profile, we counted the different arrangements of methylated and unmethylated cytosines found in sequencing reads spanning the entire epilocus. We then grouped the observed arrangements in five methylation classes (MCs) according to the number of methylated cytosines they bear. Hence, we computed the epilocus MC profile as the fraction of reads supporting a given MC out of the total number of reads (Figure 1B).

Measure of dissimilarity. We adopted the Jensen-Shannon Distance as a measure of dissimilarity between two MC profiles. The Jensen Shannon Distance (JSD) quantifies the degree of dissimilarity between discrete distributions P_1 and P_2 (42), which in our case are represented by two sets of relative abundances, and is defined as

$$d = \sqrt{\frac{D(P_1, \bar{P}) + D(P_2, \bar{P})}{2}}$$

In this formula, $\bar{P} = \frac{P_1 + P_2}{2}$ represents the average distribution of two MC profiles, which is obtained by averaging the relative abundance of each MC among the two samples. $D(P_1, \bar{P})$ and $D(P_2, \bar{P})$ represent the Kullback-Leibler (KL) divergence between the average profile \bar{P} and the profile P_1 and P_2 , respectively. The KL divergence of two discrete distributions is computed as $D(A, B) = A * \log_2(\frac{A}{B})$ (43).

Establishment of MC profiling thresholds. We used deep amplicon bisulfite sequencing (D-ABS) data to synthesize 4-CpG low coverage datasets, starting from an in-house database of D-ABS amplicons produced in previously published studies (27,44–46). Detailed descriptions of the employed amplicons can be found in Supplementary Table S1.

First, we split each amplicon into non-overlapping regions made up of four CpGs, thus obtaining several 4-CpG high-coverage datasets. Among these datasets, we selected those with higher coverage (number of reads > 20 000). Since we expect that fully methylated or unmethylated profiles would be better captured at low coverage than intermediately methylated ones, due to the higher number of methylation classes with non-zero abundance, we selected 4-CpG datasets with average methylation levels spanning the entire range from 0 to 1 and enriched for datasets with intermediate average methylation. In this way, we selected 25 datasets, representative of five groups according to the average methylation level (Supplementary Tables S2 and S3).

To simulate low coverage datasets, we randomly sampled a fixed number of reads from each 4-CpG dataset.

We adopted the low coverage datasets to address the following issues:

- 1) the minimum coverage to minimize the error between a reference MC profile (i.e. the MC profile computed from the high coverage dataset) and an estimated MC profile (i.e. the MC profile computed from a low coverage dataset)

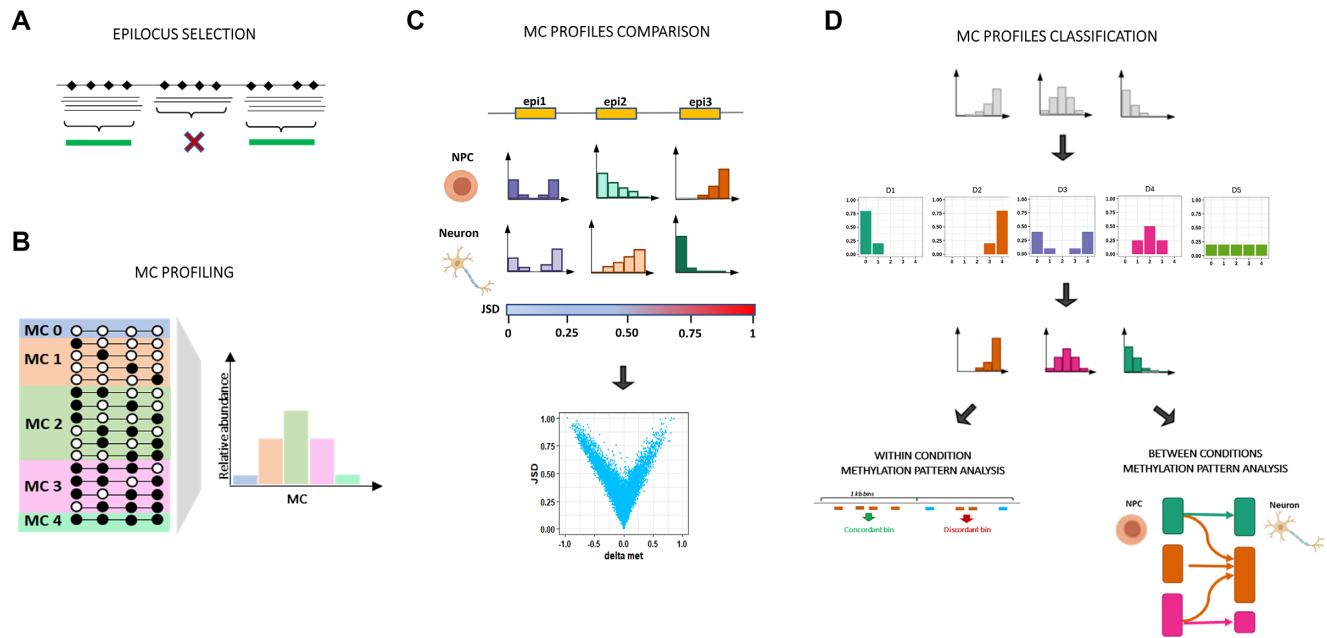


Figure 1. Schematic drawing of MC profiling. (A) Selection of epiloci eligible for MC profiling. An epilocus is defined as a genomic region holding 4 CpGs. Non-overlapping epiloci with coverage higher than 50 reads were retained for MC profiling. (B) For a given epilocus, the MC profile is computed as the fraction of reads supporting the possible MCs (i.e. groups of molecules bearing a given number of methylated cytosines, independently of the position) out of the total number of reads. Only reads spanning the entire epilocus were considered in the computation of the MC profile. (C) The Jensen-Shannon distance is used to quantify the degree of dissimilarity between MC profiles. Based on the results obtained from simulated data, we considered two MC profiles to be different when observing a JSD above 0.26. The JSD can be used to assess the changes of MC profiles at a given epilocus in different conditions. The JSD can be also compared to other metrics, such as the difference of average methylation (delta met), over the analyzed epiloci. (D) MC profiles were assigned to 5 Methylation Patterns (MPs) according to the most similar among 5 archetypal profiles (here indicated in the upper panel, middle row). This data compression procedure provided us with a signature of genome-wide MC profiles composition in a given condition. MPs enabled us to i) directly compare the MP of different epiloci within the same sample/ or condition (within sample analysis) and ii) to compare the MP transitions occurring at a given epilocus in different conditions (between conditions analysis).

To address this point, we synthesized 1000 low coverage four CpG datasets for coverage values ranging from 20 to 200. From each dataset, we calculated the MC profile, and computed the JSD from the MC profile of the respective four CpG high-coverage dataset. As shown in Supplementary Figure S1A, the JSD values decreased as the coverage increased, as expected. In particular, JSD values dropped between 25 and 50 reads (Supplementary Figure S1A). A similar gain in accuracy is achieved by triplicating the coverage (i.e. achieving a read number higher than 150). Based on these observations, we considered a region covered by at least 50 reads to be eligible for MC profiling.

- 2) the minimum value of JSD to consider 2 MC profiles as different.

To address this point, we simulated 1000 pairs of low-coverage datasets with a fixed coverage of 50 reads. For each dataset pair, we computed the JSD among the estimated MC profiles. Ideally, two read groups sampled from the same dataset should exhibit very similar, if not identical, profiles, with a JSD value approaching 0. In practice, however, the estimated profiles differed to a certain extent. As shown in Supplementary Figure S1B, at a coverage of 50 reads, MC profiles exhibited a JSD lower than 0.26 (min = 0.22, max = 0.28) for 95% of the experiments. JSD values observed for a wider range of coverage are reported in Supplementary Table S4.

We concluded that MC profiles having a JSD higher than 0.26 could be defined as different with an error equal or lower than 0.05. Hence, when comparing two MC profiles, we considered them to be different when we observed a JSD above 0.26.

Epilocus filtering and multiple samples handling. In individual samples, epiloci with coverage lower than 50 reads and higher than 99th percentile were filtered out. Overlapping epiloci were also removed (Figure 1A).

To handle MC profiles observed in different samples at a given epilocus, we first computed the JSD between the possible sample pairs and retained those epiloci with JSD below 0.26 in all the pairs. For these epiloci, we computed the average MC profile by averaging the relative abundance of each MC among the samples.

In this way, we obtained a consensus MC profile representative of all samples in a given condition. This enabled us to directly compare consensus profiles of an epilocus in two conditions through the JSD (Figure 1C).

MC profiles classification. To provide biological interpretation of MC profiles, we adopted a data compression scheme. We assigned each MC profile to one among 5 methylation patterns (MP) according to the most similar of 5 archetypal profiles (Figure 1D, see Materials and Methods). These reference profiles, named from D1 to D5, are reminiscent of standard discrete distributions and were cho-

sen because they reflect the reasonable profiles of an epilocus expected at a given methylation amount. In particular, D1 and D2 represent the two specular profiles for highly methylated or unmethylated epiloci, in which we expect a prevalence of fully unmethylated and methylated MCs, respectively. D3, D4 and D5, instead, represent the hypothetical profiles of intermediately methylated regions, that can reflect: (i) the prevalence of both fully methylated and unmethylated MCs (D3, bimodal profile), (ii) the prevalence of intermediately methylated MCs (bell-shaped profile, D4) or (iii) the presence of all possible MCs with the same relative abundance (uniform profile, D5).

The data compression procedure provided us with a signature of MC profiles composition over all the analyzed epiloci of a sample. Importantly, this signature did not depend on the specific experimental system. In this way, we were able to (i) directly compare the MP of genome-wide epiloci within the same sample/condition and (ii) to compare the MP transitions occurring at a given epilocus in different conditions (Figure 1D).

To assign an MC profile to the nearest MP, we computed its JSD from the five prototypes, and assigned it to the MP corresponding to the prototype with minimum JSD (Figure 1D). To check the appropriateness of our classification procedure, we compared the JSD of each MC profiles from the two nearest MPs, with the lower JSD value representing the distance from the membership pattern centroid (Within Class Distance, WCD) and the second value representing the distance from the nearest outer pattern centroid (External Class Distance, ECD). We observed that the ECD-WCD ratio exceeded 1.5 for 98% of MC profiles in Dataset 1 and 95% of MC profiles in Dataset 2 (Supplementary Figure S2). Thus, we concluded that the proposed scheme was roughly consistent and representative of the diverse MC profiles observed in our datasets.

Dataset

We analyzed previously published RRBS data and enhanced RRBS data (47–50). The data were publicly available in the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>) with the following accessions: GSE66121, GSE130735, GSE53714, GSE72700. When a consistent discrepancy existed for the number of epiloci covered by at least 50 reads in different samples, we retained those samples providing the highest number of epiloci. A detailed description of the samples adopted from each dataset is included in Supplementary Table S5.

Normalized FPKM expression values for human CD19 + normal B-cells were obtained from GEO with the accession GSE66121.

Data processing

RRBS raw data processing. Raw RRBS data were processed using an in-house pipeline. Fastq files were first quality checked by using FastQC (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Low-quality bases were removed using Trim Galore v0.6.6 with parameters `-rrbs` and `-paired` for paired end experiments (<https://www.bioinformatics.babraham.ac.uk/projects/trim-galore/>). The

obtained reads were aligned to the hg19 or mm10 reference genomes by using Bismark v0.23.0 with default parameters (51). The obtained BAM files were then sorted and indexed using the SAMtoolsKit (<http://www.htslib.org/>).

Deep: amplicon bisulfite sequencing data processing. D-ABS data were processed as previously described (27,44,45). In brief, paired-end reads were merged in a single fastq file by using PEAR (minimum overlapping residues equal to 40) (<https://cme.h-its.org/exelixis/web/software/pear/doc.html>). The fastq file was then converted to fasta by using PRINSEQ (<http://prinseq.sourceforge.net/>).

Epiallele counts extraction. For RRBS data, epiallele counts were extracted from BAM files using the *EpiStatProfiler* R package ((52), <https://github.com/BioinfoUninaScala/epistats>). First, genomic regions covered by at least 50 reads were individuated through the *filterByCoverage* function. Target regions holding 4 CpGs (the epiloci described in this manuscript), stepping up to 1 CpG at time, were then defined by using the *makeBins* function. The maximum length of the target regions was set from 70 to 100 bp, depending on the specific library design. Epiloci covered by at least 50 reads spanning the entire region were retained. Finally, selected epiloci were analyzed by using the *epiStatAnalysis* function with default parameters. For each epilocus, the function returns a table with summary statistics (including the average methylation), and a file with epiallele counts. This latter was then analyzed through in-house R scripts to compute the epilocus MC profile, as described above.

For D-ABS data, epiallele counts were then extracted by using the *AmpliMethProfiler* tool (53). The MC profile of the amplicon was then computed following the same procedure of RRBS epiloci.

Allele-specific alignment sorting. To perform allele specific MC profiling, we applied the pipeline based on the *SNPsplit* tool (54) on a dataset of crossed strain mice. First, the positions holding alternative sequences between the strains were extracted from the VCF file downloaded from the Mouse Genomes Project repository (ftp://ftp-mouse.sanger.ac.uk/current.snps/mgp.v5.merged.snps_all.dbSNP142.vcf.gz), and were masked from the reference mm10 genome by using the *SNPsplit.genome.preparation* function in single strain mode. Fastq files were then aligned to the masked genome by using Bismark 0.23.0 with default parameters. The reads aligned to polymorphic sites were assigned to the respective allele by using the *SNPsplit* function. In brief, the reads aligned to variant positions were tagged (*SNPsplit-tag* internal function), assigned to the reference or to the alternative allele (*tag2sort* internal function), and written down in separate bam files. We ran the *SNPsplit* function in `-bisulfite` mode to automatically discard the reads aligned to C/T or T/C variants on the forward strand and to G/A or A/G variants on the reverse strand, since these variants cannot be distinguished from a methylation status call. The bam files relative to the reference and the alternative allele were processed independently with the *EpiStatProfiler* tool to obtain the epiallele counts and to compute the MC profile. At the end, we were able to profile

2749, 460, 314 autosomal epiloci in three mice, with a minimum coverage of 50 reads on both alleles.

Epiloci annotation. Epiloci were annotated by using the *annotatr* R package against hg19 and mm10 CpG tracks and hg19 and mm10 genes tracks.

Epiloci were associated with the nearest genes by using the *seq2pathway* R package. To minimize the number of genes associated with an epilocus, the ‘adjacent’ parameter was adopted, thus enabling to assign each epilocus to the closest genes only. For the association, the *FullResult* output was considered, which also included non-coding genes.

To test the association between MC profiles and chromatin marks changes, epiloci were annotated using the chromHMM segmentation tracks produced for the GM12878 lymphoblastoid cell line from the RoadMap Epigenomics project (<https://egg2.wustl.edu/roadmap/data/byFileType/chromhmmSegmentations/ChmmModels/coreMarks/jointModel/final>), whereas epiloci of Dataset 5 were annotated using the segmentation tracks for mouse hindbrain (E10 and P0) downloaded from UCSC (55). Each epilocus was annotated with the label of the genomic segment with the highest overlap. The E10 hindbrain track was used to annotate epiloci of hippocampal precursors, whereas the P0 hindbrain track was adopted to annotate epiloci in Granule cells and CA neurons.

Association of MPs with expression level. We used the normalized expression data (FPKM) available for 3 samples from Dataset 2 to assign genes with expression categories. For each gene, we computed the average value among the samples. We then labeled as highly-expressed genes those with expression value above the median, and labeled as lowly expressed genes the ones below or equal to the median.

Epiloci were assigned to gene promoters, exonic or intronic regions by using the *annotatr* R package against hg19 genes track.

The association between the number of epiloci assigned to the different MPs and the expression status was tested through chi-square test and post-hoc analysis of chi-square residuals (see Statistical test).

Statistical analysis

Classification concordance of neighboring epiloci. To test the classification concordance of neighboring epiloci, we first binned the genome into 1 kb long regions. We then intersected the bins’ coordinates with that of epiloci shared by all the samples in the dataset (see Epilocus filtering and multiple samples handling). We removed the bins harboring less than three epiloci and labeled the remaining ones as concordant if they hold epiloci assigned to the same prototype class, and discordant otherwise. We tested the hypothesis that the number of concordant bins was higher than the one expected by chance by bootstrapping. In brief, we scrambled the epiloci grouped in each bin, such that the overall number of bins together with the number of epiloci they hold reflected those observed in experimental data, but the epiloci were no longer grouped in a bin based on their

proximity but were randomly sampled without replacement from the dataset. We repeated this procedure 1000 times, and each time we counted the number of scrambled bins classified as concordant. We thus obtained the distribution of the number of concordant bins expected by chance, that we compared with the number of concordant bins observed in experimental data.

MC profiles heterogeneity. For haploid models, we adopted the epilocus MC counts, i.e. the number of MCs with non-zero relative abundance, to estimate the degree of cellular heterogeneity of DNA methylation. The MC counts distribution of haploid epiloci was compared to that of diploid epiloci to estimate the contribution of allelic heterogeneity to MC profiles. For male X epiloci, autosomal epiloci in the same sample were used as diploid reference, whereas for polymorphic epiloci the joined MC profiles were used. When directly comparing the MC profiles of an epilocus on the two alleles in the second model, we controlled for coverage differences. We found no significant differences between the two alleles (paired Wilcoxon test P -values < 0.01).

Statistical test. All the statistical analyses were performed using R software (version 4.0) with an alpha value set for $P < 0.01$.

Association between categorical variables was tested for statistical significance through Fisher test (when both categorical variables were dichotomous). We applied it to test whether MC profile changes more probably involved epiloci that also underwent chromatin changes upon differentiation, epiloci located in promoters or epiloci located in CpG Islands.

Association between non-dichotomous categorical variables was tested for statistical significance through chi-square test and post-hoc analysis of chi-square residuals (*chi.square.posthoc.test* function from the homonymous R package, adopting Bonferroni correction to control for alpha inflation). We applied chi-square to test whether epiloci exhibiting inter-individual variability were enriched in peculiar genomic contexts (promoters, exons, introns, or intergenic regions), or whether epiloci assigned to different MPs were enriched in particular genomic regions (for example, regions flanking imprinted genes or regions decorated with different histone marks) or more probably changed MC profiles upon differentiation.

Differences in average reciprocal distance among epiloci in concordant and discordant bins was tested through the Mann-Whitney test.

Enrichment analysis for 5129 epiloci with significant changes in MC profiles and stable average methylation upon differentiation was performed using GREAT version 4.0.4 (56), using default parameters and the coordinates of all the analyzed epiloci (115 608) as background.

RESULTS

The MC profiling approach

Rationale of MC profiling. The rationale of MC profiling, and the differences with epiallele-based approaches, is depicted in Figure 2.

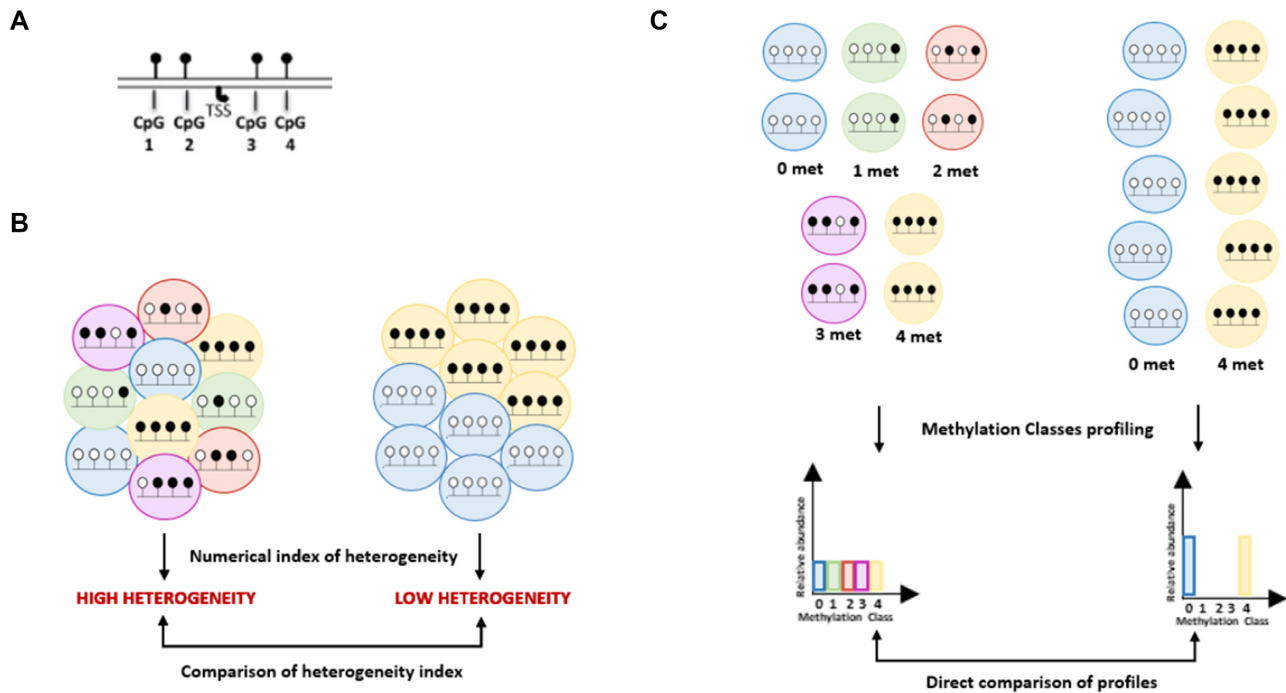


Figure 2. Rationale of MC profiling. (A) Example of a region of interest holding four CpGs (TSS = Transcription Starting Site) (B) Representation of epiallele-based analysis. The DNA methylation heterogeneity for a certain locus is usually quantified through a numerical index (e.g. epipolymorphism, Shannon entropy, ...). This can be then adopted to compare the heterogeneity of pools of molecules (for example, to compare the heterogeneity of a certain locus in different samples). (C) Representation of MC profiling analysis. The epialleles are first grouped in methylation classes (MCs) according to the number of methylated cytosines. The relative abundances of the possible MCs (for a locus holding n CpGs there are $n + 1$ possible MCs), named MC profile, summarize the molecular heterogeneity and the methylation levels of a given region. MC profiles can be directly adopted to perform differential analysis.

Epiallele-based approaches are based on the direct analysis of the arrangements of methylated and unmethylated cytosines (epialleles) in sequencing reads mapped to a region of interest (Figure 2A).

Considering each reads coming from a DNA molecule, several scores have been developed to quantify the heterogeneity observed in a bulk sample, and to compare it among different samples (Figure 2B) (19). This approach has proved to be particularly suitable, for example, to individuate regions undergoing clonal selection and epigenetic drift in tumors (33,34,57). In this setting, the composition of individual epialleles is only indirectly accounted for. Similar heterogeneity values could, indeed, come from different epiallele compositions. Of note, the methylation level of epialleles is usually not, or only partially, incorporated in these heterogeneity scores, which makes difficult to interpret the functional impact of heterogeneity shifts (33–36,57).

The underlying idea of our approach is that looking at the distribution of epialleles grouped by their methylation levels adds useful information for the functional interpretation of DNA methylation heterogeneity in a sample. The proposed approach, MC profiling, is indeed based on the empirical estimate of the distribution of epialleles grouped by their methylation levels (Figure 2C). We already applied the concept of MCs in previous works with the aim to model DNA methylation dynamics at targeted loci assayed through high-coverage bisulfite sequencing (40,41). We here extended our approach to enrichment-based genome-wide datasets, like

the ones from Reduced Representation Bisulfite Sequencing (RRBS) experiments, thus allowing for the simultaneous analysis of thousands of regions from the same sample. In this context, we implemented a new analytical framework to directly compare MC profiles across regions and samples.

Instead of adopting a numerical index (as, for example, the Shannon Index) to summarize the DNA methylation heterogeneity of a given region, we kept as much information as possible and described, for each DNA region, the relative abundance of the possible MCs. In this setting, we adopted the direct comparison of MC profiles to analyze differential methylation of a given region among conditions, or to examine the differences among regions in the same condition (Figure 2C). It is important to point out that, in this latter setting, direct comparison of epiallele composition would only be possible through MCs, being these sequence independent, and not through the epialleles themselves.

Comparing MC profiles allowed us indeed to compare not only the heterogeneity but also the different methylation levels of DNA molecules.

In summary, adopting MC profiles can provide the following advantages:

- Considering how they are computed, MC profiles directly incorporate the average methylation of a given region, and inform on how it is contributed by single DNA molecules.

- MC profiles retain all information from a pool of molecules, and enable the direct visualization of DNA methylation heterogeneity of a given region
- MC profiles are empirically estimated from sequencing reads, and are independent on a priori parametrization of DNA methylation dynamics (see Discussion)

Description. In this study, we focused on the methylation status of regions made up of four CpG sites, moving from the observation that the number of reads per region drops when increasing the number of CpG sites from 4 to 5 (34). We will refer to these regions as epiloci in the following text (Figure 1A). Based on the results obtained on simulated data (see Methods), we selected for MC profiling those epiloci covered by at least 50 reads, considering only reads spanning the entire epilocus.

For a certain epilocus, methylated and unmethylated cytosines can be arranged in 16 possible combinations which can be in turn grouped in methylation classes (MCs) according to the number of methylated cytosines they bear. For a four CpG-locus, five MCs can be indeed described. We depicted the DNA methylation status of a given epilocus through its inherent MC profile, i.e. the set of the relative abundances of the five MCs. To compute this MC profile, we counted the different arrangements of methylated and unmethylated cytosines found in sequencing reads spanning the entire epilocus (Figure 1A, B).

Throughout this study, we adopted the Jensen Shannon Distance (JSD) to quantify the dissimilarity between MC profiles (see Materials and Methods).

Based on the results of simulations performed on high-coverage targeted bisulfite sequencing data, we considered eligible for MC profiling those epiloci covered by at least 50 reads, and we considered two profiles to be different when we observed a JSD above 0.26 (see Methods).

To improve the interpretability of the data, we adopted a data compression procedure, and assigned each MC profile to a Methylation Pattern (MP) according to the most similar of 5 archetypal profiles (Figure 1D), hereafter referred to as prototypes (see Methods). The prototypes, which are reminiscent of standard discrete distributions, were chosen because they reflect the reasonable profiles of an epilocus expected at a given methylation amount. In fact, D1 and D2 represent the two symmetric profiles for highly methylated or unmethylated epiloci, in which we expect a prevalence of fully unmethylated and methylated MCs, respectively. D3, D4 and D5, instead, represent the hypothetical profiles of intermediately methylated regions, that can reflect (i) the prevalence of both fully methylated and unmethylated MCs (D3, bimodal profile), (ii) the prevalence of intermediately methylated MCs (bell-shaped profile, D4) or (iii) the presence of all possible MCs with the same relative abundance (uniform profile, D5).

MC profiles conjugate quantitative methylation and molecular heterogeneity of an epilocus. We applied MC profiling to two datasets of samples publicly available in GEO (see Materials and Supplementary Table S5). Dataset1 included samples from 3 wild-type mice embryos, whereas Dataset2 included three samples from human CD19+ B-cells isolated from normal controls. Indeed, our datasets came from dif-

ferent species and were representative of different developmental stages, where we expect that DNA methylation heterogeneity probably derives from different dynamics (epigenetic drift in somatic cells versus cell differentiation in mouse embryos). We reasoned that such an experimental plan would have enabled us to generalize the results of our analysis.

For each sample, we profiled about 100 000 epiloci in Dataset 1 and 90 000 epiloci in Dataset 2. The systematic description of the analyzed cytosines is reported in Supplementary Figure S3 and S4.

By examining the average methylation of epiloci belonging to different MPs, we confirmed that the quantitative amount of methylated cytosines of assigned elements was coherent with the expected values for each pattern (Figure 3A, Supplementary Figure S5). However, MC profiles add further information depicting the heterogeneity of DNA methylation among DNA molecules. This was particularly evident for the D3, D4 and D5 patterns. In fact, epiloci exhibiting the same average methylation were assigned to different MPs (Figure 3B).

MC profiles are mostly stable among individuals and across genomic regions. We investigated the stability of MC profiles across samples. For this aim, we analyzed the epiloci for which the MC profiles were assessed in all the samples in the individual datasets ($n = 87\,457$ and $n = 41\,609$) and computed the JSD of MC profiles among sample pairs. We found that 98% of epiloci in Dataset 1 and 96% in Dataset 2 had JSD lower or equal to 0.26 in all sample pairs, meaning that MC profiles at most of the epiloci were very similar between samples (Figure 3C, Supplementary Figure S6A).

In both datasets, we found that stable epiloci, i.e. epiloci with a JSD below the cutoff in all sample pairs ($n = 86\,319$ and $n = 39\,767$, in Datasets 1 and 2 respectively), were enriched in promoters (chi-square post hoc test P -values $< 1e-7$) and depleted in intergenic regions (chi-square post hoc test P -values $< 1e-7$). On the contrary, variant epiloci, i.e. epiloci with JSD above the cutoff in at least one sample pair ($n = 1138$ and $n = 1842$, in Datasets 1 and 2 respectively), were depleted in promoters (chi-square post hoc test P -values $< 1e-7$) and were enriched in intergenic regions (chi-square post hoc test P -values $< 1e-7$). We found no difference in the proportion of stable and variant epiloci located in coding sequences (Figure 3D, Supplementary Figure S6B).

Based on this result, for each dataset we retained for further analysis the stable epiloci, and computed the consensus MC profile by averaging the relative abundances of each MC from the three samples. We then applied the data compression procedure, and assigned the consensus MC profiles to the MPs (see Materials and Methods).

Since epigenetic modifications are expected to involve larger DNA regions than individual epiloci, we expected that neighboring epiloci exhibited concordant MC profiles. To test this hypothesis, we binned the genome in 1 kb regions, and compared the MPs of epiloci located in each bin (see Materials and Methods). Among the bins harboring at least three epiloci, 10 733 (99%) bore concordant and (1%) 148 bore discordant epiloci in Dataset 1, whereas (95%) 5079 bore concordant and (5%) 249 bore discordant epiloci

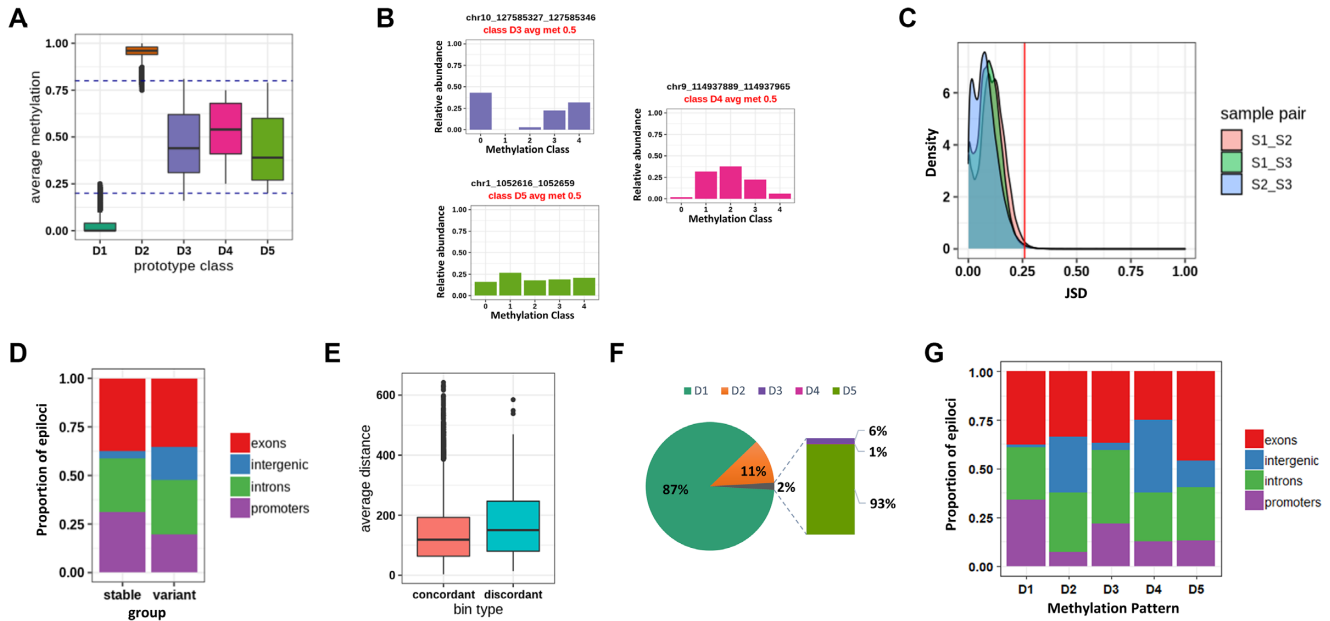


Figure 3. MC profiling results for Dataset 1. (A) Average methylation level of epiloci assigned to each MP. (B) Example of epiloci with same average methylation and different MC profiles. (C) Density plot of MC profile distance between sample pairs. x-axis: JSD values between sample pairs. y-axis: density of epiloci with a given sample-pairs JSD value. The red line indicates the cutoff value of JSD. (D) Genomic annotation of epiloci with stable or variant MC profiles. (E) Average distance between epiloci inside concordant and discordant bins. (F) Fraction of epiloci attributed to the different MPs. (G) Genomic annotation of epiloci assigned to the different MPs.

Dataset 2. We confirmed that the number of bins bearing concordant epiloci significantly differed from the one expected by chance in both datasets (see Materials and Methods and Supplementary Figure S7). This result suggests that MC profiles of neighboring epiloci tend to be similar. This conclusion is further supported by the observation that the reciprocal distance between epiloci in concordant bins tends to be lower than in discordant bins (Mann-Whitney P -value = 0.0005866, Figure 3E, Supplementary Figure S6C).

Overall, MC profiles resulted to be mostly stable among individuals and across genomic regions, thus suggesting that the heterogeneity captured by MC profiles mostly results from controlled DNA methylation dynamics, rather than from stochastic fluctuations of methylation levels.

MC profiles differentiate functional genomic regions. We reasoned that assigning MC profiles to different MPs could provide us with a signature of genome-wide MC profiles composition in a given dataset. We indeed examined the proportion of epiloci assigned to each MP. In accordance with the well-established bimodal distribution of average DNA methylation (16), the most represented prototype classes were D1 (83% and 78% of epiloci in Datasets 1 and 2, respectively) and D2 (about 15% and 16% of epiloci in Datasets 1 and 2, respectively). The intermediately methylated D3, D4 and D5 classes accounted respectively for 2% of epiloci in Dataset 1 and 5% of epiloci in Dataset 2 (Figure 3F, Supplementary Figure S6D). Among the intermediately methylated classes, the most represented one was the D5 (90% and 82% of epiloci in Datasets 1 and 2, respectively), followed by the D3 class (9% and 13% of epiloci in

Datasets 1 and 2, respectively). The D4 class was strongly underrepresented (1% and 4% of intermediately methylated epiloci in Dataset 1 and 2, respectively) in normal conditions (Figure 3F, Supplementary Figure S6D), suggesting that intermediate values of average methylation rarely reflect an intermediate methylation amount on different DNA molecules. Instead, intermediate values of average methylation more often reflected the coexistence of fully unmethylated and fully methylated molecules, in presence (D5) or in absence (D3) of intermediately methylated molecules.

The classification of MC profiles to MPs also enabled us to investigate whether epiloci attributed to the different prototype classes were located in genomic regions with different functional characteristics. As shown in Figure 3G and Supplementary Figure S6E, we found that the D1 class was enriched within promoters and exons (chi-square post hoc P -values $< 1e-7$) and depleted in intergenic regions and introns (chi-square post hoc P -values $< 1e-7$). On the contrary, the D2 class was mainly located in intergenic regions and introns (chi-square post hoc P -value $< 1e-7$) and depleted in promoters and exons (chi-square post hoc P -value $< 1e-7$). Similarly, the D5 class was depleted from promoters and enriched in intergenic regions (chi-square post hoc P -values $< 1e-7$). We did not find significant differences in the localization of D3 and D4 epiloci.

We found that MPs composition could further distinguish genomic regions decorated with different histone marks in Dataset 2 (Supplementary Figure S8A). For example, MPs separated constitutive heterochromatin from Polycomb-repressed regions (chi-square P -value $< 1e-7$), with the former enriched not only for the methylated D2 but

also for the D5 MP (chi-square post hoc P -value $< 1e-7$), pointing to higher heterogeneity in constitutively inactive genomic regions. Polycomb-repressed regions, on the other hand, were enriched for the D1 MP (chi-square post hoc P -value $< 1e-7$), pointing to lower levels of DNA methylation in Polycomb-regulated regions.

We also found that MPs composition varied when separately investigating the promoter, exonic and intronic regions of genes with expression levels lower or higher than the median value in Dataset 2 (chi-square P -value $< 1e-7$, Supplementary Figure S8B, C). We found that D1 MP was enriched in promoters, introns and exons of highly-expressed genes (chi-square post-hoc P -value $< 1e-7$), whereas the D2 MP was enriched in exons (chi-square post-hoc P -value $< 1e-7$) and slightly enriched in promoters of lowly-expressed genes (chi-square post-hoc P -value $< 5e-3$). Again, we found an enrichment of the D5 MP in promoters, exons and introns of lowly expressed genes (chi-square post-hoc P -value $< 1e-7$), suggesting a consistent pattern of increased heterogeneity in low-to-inactive regions.

Cellular heterogeneity is the strongest contributor to MC profiles

MC profiles recapitulate heterogeneous methylation status among DNA molecules. This heterogeneity can, in principle, reflect both allelic and cellular differences. For most epiloci, these two components cannot be distinguished in a bulk experiment, in which the information on how DNA molecules are paired in individual cells is missing. We reasoned that epiloci present as single copies in the genome could be a good model to investigate the contribution of cellular differences to MC profiles. In fact, at these loci, the presence of multiple MCs can reflect only cellular differences.

As a first model of DNA regions present as single copies in the genome, we investigated epiloci located on the X chromosome of a male mouse from Dataset 1 ($n = 1303$). For each epilocus, we quantified the cellular heterogeneity in terms of MC counts, i.e. the number of MCs supported by at least 1 DNA molecule. We found MC counts above 1 for most epiloci ($>70\%$), pointing to cellular heterogeneity of DNA methylation as the rule for most epiloci. Of note, the distribution of MC counts' values for X epiloci did resemble that of autosomal epiloci. This observation seems to suggest that the degree of heterogeneity captured by MC profiles is poorly affected by the copy number of the given epilocus (Figure 4A).

As an additional model, we studied autosomal epiloci of mice born from two different strains (Dataset 3 in Supplementary Table S5). Based on known polymorphic sites between the two strains, we were able to attribute each read to the respective allele, and to explore the allele specific MC profile for more than 300 autosomal epiloci in three mice (see Materials and Methods). When analyzing the joint MC profiles, we confirmed the high degree of DNA methylation heterogeneity, with about 95% of the autosomal epiloci having an MC count equal or greater than 1 (Supplementary Figure S9A–C). However, when directly comparing the MC profiles of the reference and alternative alleles, we found no

differences for most of the epiloci, with only a small proportion of epiloci (6–7%) exhibiting allele specific methylation (Figure 4B).

Overall, the results from the analysis of male X chromosome MC profiles and allele specific MC profiles provided evidence for cell-to-cell differences as the major contributor to MC profiles, with evidence of allelic differences only in a small fraction of autosomal epiloci. In addition, we found that the allelic MC profiles of 144 epiloci shared among the samples were mostly stable among sample pairs, thus suggesting that similar patterns of cellular heterogeneity were present in different individuals (data not shown).

MC profiling individuates a signature of imprinted regions

We tested the capability of MC profiling to discriminate regions undergoing genomic imprinting, a well-known phenomenon of allele specific regulation. In these regions, it is expected that the two alleles differ for their DNA methylation status. Hence, we wondered whether D3 epiloci, in which two pools of molecules exist with opposite DNA methylation status, were enriched at genomic regions flanking imprinted genes.

To test this hypothesis, we assigned each epilocus of Datasets 1 and 2 to its nearest gene (see Methods) and marked the epiloci as associated with imprinted genes if the closest gene was enlisted in Geneimprint (<https://www.geneimprint.com/>). As shown in Figure 5A and Supplementary Figure S10, the five MPs were differentially represented among epiloci flanking imprinted and not imprinted genes (chi-square test P -values $< 2.2e-16$). Specifically, epiloci assigned to the D3 pattern were strongly overrepresented among epiloci flanking imprinted genes (chi square post-hoc test P -values $< 2e-16$), thus confirming that D3 epiloci were preferentially, even though not exclusively, associated with allele specific methylation.

As a confirmatory experiment, we searched for D3 epiloci flanking imprinted genes in Dataset 3. We found a single epilocus with these characteristics, located on chr1, upstream of the *Zdbf2* imprinted gene. In this locus, differentially methylated regions had been previously described (58). Figure 5B shows how the bimodal joint MC profile at this epilocus results from different profiles on the two alleles, with one skewed towards complete demethylation and the other towards complete methylation. It is worth noting that, for both alleles, MC profiling individuated a certain degree of cellular heterogeneity, since intermediate MCs were also represented.

Loss of genomic imprinting is a well-known epigenetic modification occurring in many tumors (59,60). Considering the association that we found between bimodal D3 MC profiles and imprinted regions, we wondered whether changes in MC profiles of D3 epiloci could be identified in tumor samples. When inspecting epiloci flanking imprinted genes in tumor samples (Dataset 4), we identified an epilocus, located at chr20:57415288–57415313, whose MC profile profoundly changed in one of the samples in respect to controls from Dataset 2. This epilocus was located in the promoter of the *GNAS* gene, for which loss of imprinting in tumors has been described (Figure 5C) (60).

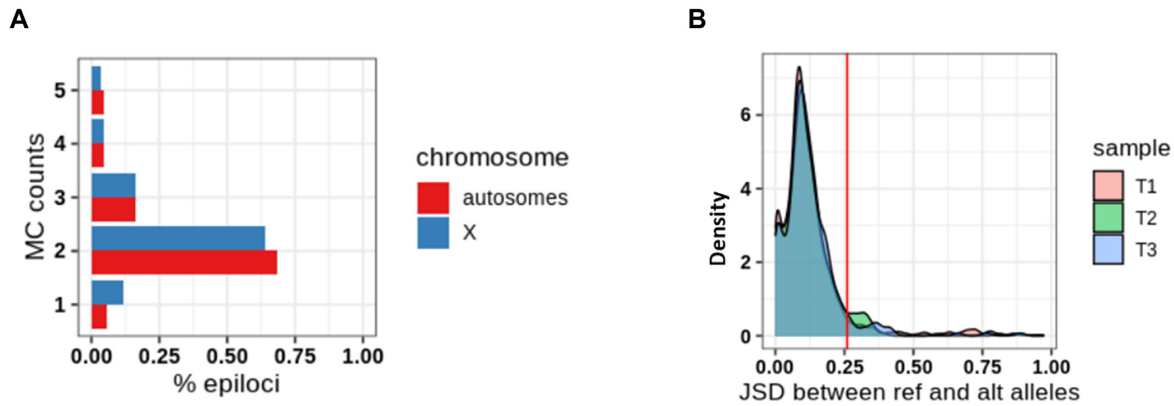


Figure 4. MC profiling of haploidy epiloci. (A) Fraction of epiloci (x-axis) exhibiting a given value of MC count (y-axis) on the X chromosome (blue) and autosomes (red) in a male sample from Dataset 1. (B) Distribution of JSD values between reference and alternative alleles in three samples from Dataset 3.

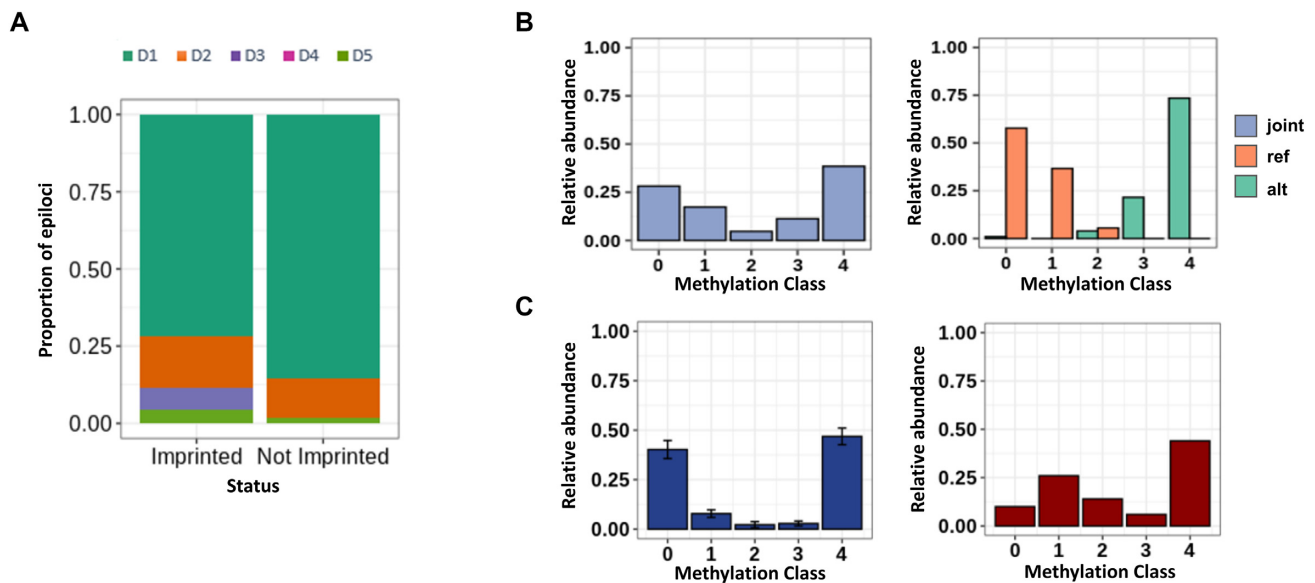


Figure 5. MC profiling of imprinted regions. (A) Proportion of epiloci assigned to the different MPs in imprinted and non imprinted genomic regions. (B) Example of bimodal epilocus flanking the *Zdbf2* imprinted gene. The joint MC profile (i.e. the profile obtained without splitting the alleles) is shown in light blue, whereas the profiles of the reference (ref) and the alternative (alt) alleles are shown in orange and green respectively. (C) Epilocus in *GNAS* promoter (chr20:57415288–57415313) with altered MC profile in a tumor sample. For this epilocus, the MC profile averaged on three control samples from Dataset 2 is shown in blue, and the MC profile from the tumor sample from Dataset 4 is shown in red.

MC profiling aids to dissect cell-to-cell differences in DNA methylation on the inactive X

Based on the results obtained from MC profiling of imprinted genes, we decided to investigate whether epiloci located on the X chromosome also exhibited peculiar MC profiles due to the X inactivation process. It is in fact known that, during the inactivation of the X chromosome, most loci are inactivated (subject loci) while others partially or totally escape this inactivation (escapee or variable escapee loci) (61).

We indeed analyzed the MPs to epiloci flanking genes with different inactivation status. First, we assigned X epiloci to the respective MP in two female samples from Dataset 2. Then, we assigned to each epilocus the consensus inactivation status of the nearest gene (61). In this way, we classified 551 epiloci as subject to X chromosome inac-

tivation, 138 as escapee, 56 as variable escapee and 233 as unknown/discordant. As shown in Figure 6A, MPs were represented in different proportions among subject, escapee and variable escapee epiloci (chi square post-hoc test P -value $< 1e-7$). Escape epiloci mostly exhibited unmethylated D1 profiles (chi square post-hoc test P -value $< 1e-7$), whereas subject epiloci mostly exhibited either bimodal D3 or uniform D5 profiles (chi square test post-hoc P -value $< 1e-7$). Both groups of MPs (D1 and D3/D5) were represented among variable escapee epiloci, none of them significantly enriched.

We hence decided to investigate the MC profile of the inactive X in Dataset 1, for which two female and one male sample was available. We reasoned that we could deduce the profile of the female inactive X by comparing the MC profile of X epiloci in males and females, and that such deduc-

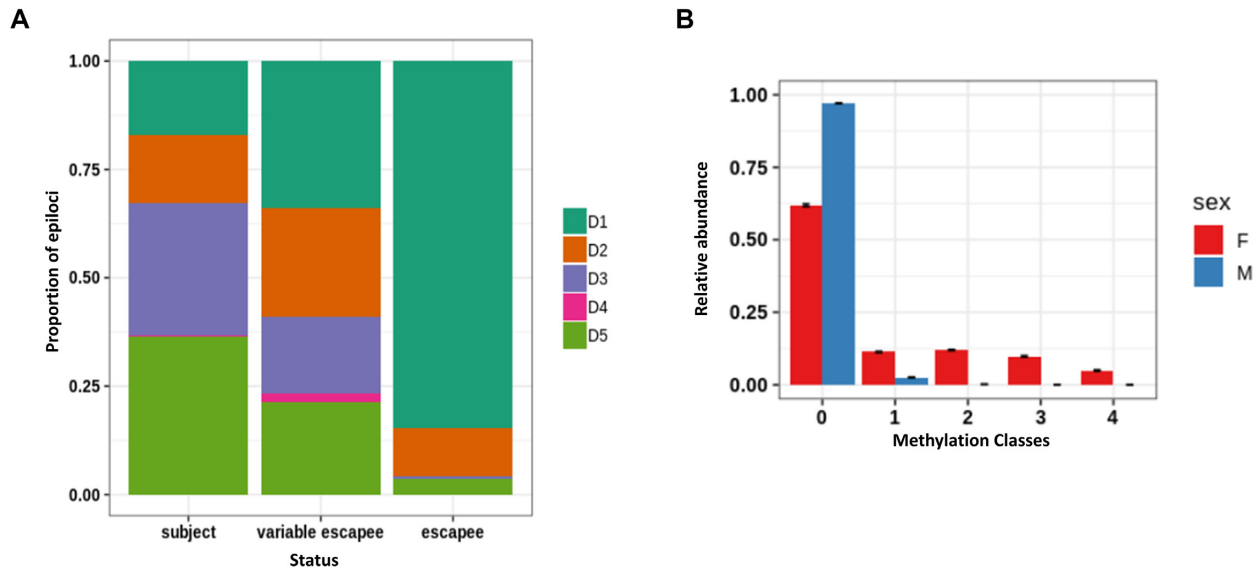


Figure 6. MC profiling of X chromosome epiloci. (A) Classification of epiloci flanking genes undergoing X inactivation (subject), stably escaping X inactivation (escapee), or variably escaping X inactivation (variable escapee). (B) Average MC profile of epiloci classified as D1 in a male sample from Dataset 1 (blue) compared to the average profile of the same epiloci in two female samples from the same dataset (red).

tion would have been particularly feasible for epiloci classified as D1 in the male samples. In fact, in this condition, it could be reasonably inferred that methylated molecules in females mostly resemble the methylation status of the inactive X. We indeed compared the average profiles of 1068 epiloci belonging to the D1 MP in males with the respective average profile in female samples (Figure 6B). The sex difference among the average MC profiles pointed to a heterogeneous DNA methylation status of the inactive X, ranging from being lowly to fully methylated in different cells. Of note, we observed a more gradual methylation status of the inactive X compared to the methylated alleles of imprinted epiloci. This observation is compatible with the previously described discrepancy of average methylation between imprinted and X inactivated genes. In fact, while for imprinted loci one of the alleles is fully methylated, X inactivated genes exhibit partial methylation of the inactive allele (58). In addition, MC profiles suggest that this partial methylation is due to cell-to-cell differences, and not to a partial methylation in all cells.

MC profiling individuates loci undergoing epigenetic remodeling upon neuronal differentiation

We challenged the ability of MC profiling to capture epigenetic changes among conditions. As a model of epigenetic changes, we choose a dataset of neuronal differentiation.

To this aim, we analyzed MC profiles changes of 115608 epiloci upon differentiation of hippocampal precursors (HP) to granule cells (GC) (Dataset 5). For each epilocus, we calculated the difference of average methylation between differentiated cells and neuronal precursors (delta meth), and quantified the MC profiles' change by using the Jensen-Shannon distance (JSD). The relationship between these two measures is shown in Figure 7A. The red lines delineate the difference of average methylation observed in 95%

of the considered epiloci (0.14), and the black line indicates the JSD threshold (0.26).

As expected, MC profiles' and average methylation changes were mostly correlated. This relationship strengthened as differences in average methylation approached the maximum, consistent with the fact that huge differences in the amount of methylated cytosines are expected to affect both average methylation and MC profiles. Symmetrically, the relationship between average methylation and MC profiles' changes weakened for lower values of average DNA methylation and was almost lost below 0.14. In this range, despite a large number of epiloci exhibiting stable MC profiles ($n = 104\,720$), a group of 5129 epiloci exhibited significant changes in MC profiles upon differentiation (blue dots in Figure 7A). As examples, Figure 7B shows two epiloci with significant changes of MC profiles and little variation of average DNA methylation. This result suggests that, at these epiloci, MC profiles were remodeled without a significant gain or loss of overall DNA methylation. The results of the enrichment analysis for genes flanking this group of epiloci is shown in Supplementary Figure S11.

To test the association between changes in MC profiles and the process of neuronal differentiation, we checked the consistency of the MC profiles changes upon differentiation of the same precursor in a different type of neuron. Notably, we found a high correlation between MC profiles changes for the 97119 epiloci examined upon differentiation of hippocampal precursors to granule cells or CA neurons (Pearson R 0.81, Figure 7C), according to the previously described high similarity among these differentiation processes (50).

To further establish the relationship between the changes in MC profiles and epigenetic remodeling upon cell differentiation, we explored the chromatin landscape, summarized by chromHMM labels, associated with the analyzed

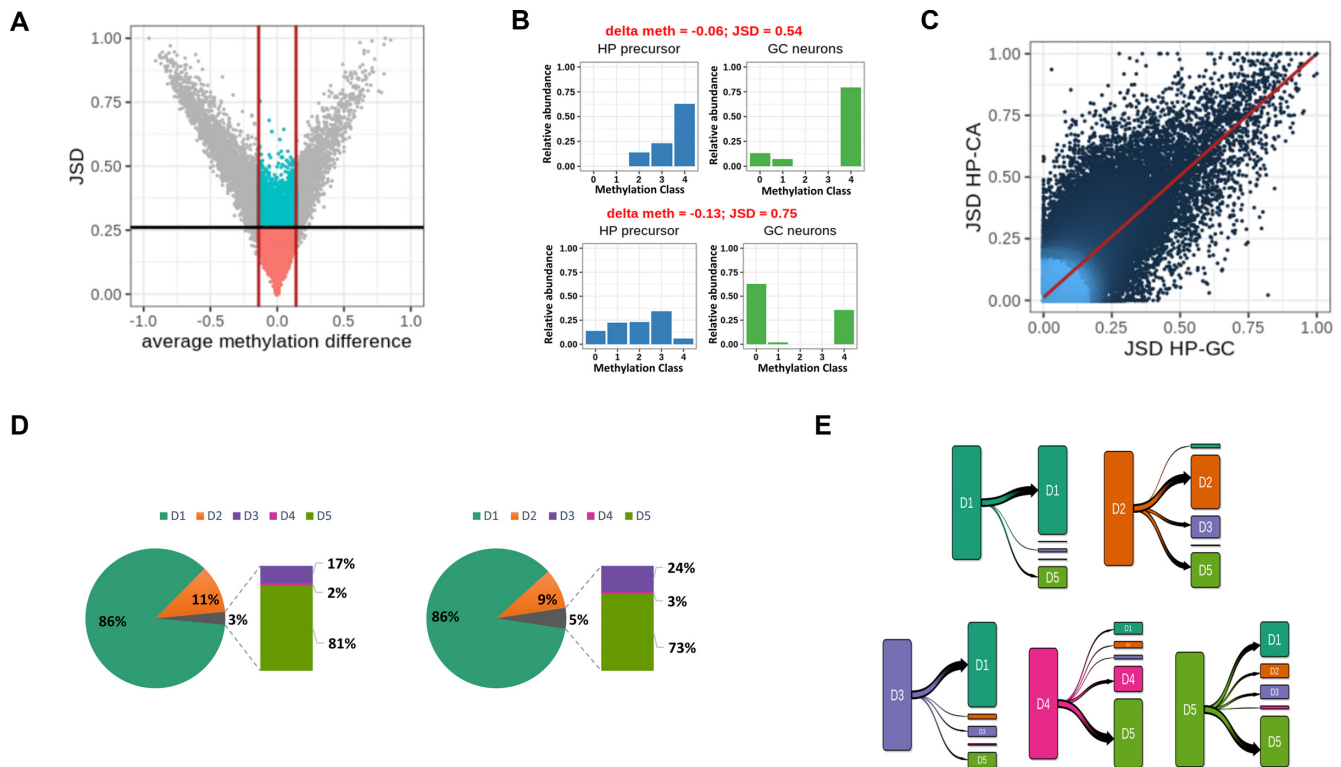


Figure 7. Application of MC profiling on neuronal differentiation. (A) DNA methylation changes at 104720 epiloci upon differentiation of hippocampal precursors to granule cells. X-axis: average methylation change (delta met); y-axis: MC profile change (JSD). The black line indicates the JSD cutoff, whereas the red lines indicate the 95th percentile of observed delta values. (B) Examples of epiloci with low difference in average methylation but high JSD values between HP and GC MC profiles (epiloci coordinates: chr19:57700749–57700788 and chr1:186924297–186924337). (C) Comparison of MC profile changes upon differentiation of hippocampal precursors (HP) to granule cells (GC) or CA3 neurons. x-axis: JSD values between MC profiles in HP and GC. y-axis: JSD values between MC profiles in HP and CA. (D) MPs composition of hippocampal precursors (on the left) and granule cells (on the right) samples. (E) Transition plot of variant epiloci in the HP-GC pair. For the epiloci assigned to the different MPs in HP cells the classification in differentiated GC neurons is shown.

epiloci. We observed that MC profile changes more probably involved epiloci located in regions that also underwent chromatin changes upon differentiation (Fisher test P -value $< 2.2 \times 10^{-16}$). In fact, only 5% of epiloci located in genomic regions with stable chromatin marks underwent changes in their MC profile, whereas 22% of epiloci undergoing chromatin changes also changed their MC profile, thus suggesting that our approach was probably identifying loci undergoing epigenetic remodeling.

When investigating the genomic localization of developmentally variant epiloci discovered by our approach, we found that they were slightly depleted outside CpG islands and promoters (Fisher test P -values $< 2.2 \times 10^{-16}$) both in HP-GC and HP-CA transitions.

Being JSD is a symmetric distance, it only quantifies the dissimilarity between two MC profiles, but does not return the information on whether this dissimilarity corresponds to a gain or loss of DNA methylation. Thus, we turned to the analysis of prototype classes to qualitatively interpret MC profile changes upon differentiation. The prototype class composition for HP and GC is shown in Figure 7D.

First, we asked whether changes were occurring at epiloci exhibiting peculiar MC profiles in neural precursors. We found that epiloci classified as D1 remained mostly stable,

whereas epiloci assigned to the other classes mostly changed their MC profile upon differentiation (chi-square post-hoc P -values $< 1 \times 10^{-7}$).

We then analyzed the prototype class composition in differentiated neurons, and found a depletion of D2 epiloci and an increased fraction of D5 epiloci (chi-square post-hoc P -values $< 1 \times 10^{-7}$), suggesting that the methylated status in differentiated neurons tends to be more heterogeneous among different cells.

Finally, to better characterize how MC profiles changes were occurring, we analyzed the prototype class transitions upon differentiation. In Figure 7E, for each prototype class in neuronal precursors, we show the final prototype class in differentiated neurons.

We noticed that for a consistent fraction of D1 and D2 epiloci, MC profiles' changes did not correspond to class transitions, meaning that these epiloci were shifting toward a higher or lower DNA methylation heterogeneity. We also noticed that a reduced fraction of epiloci evolved toward the D2 class upon differentiation.

Interestingly, most of D3 epiloci evolved to lower methylation upon differentiation, transiting to the D1 class. Thanks to prototype class analysis, we could interpret this demethylation as a negative selection of the fully methylated molecules that were present in neural precursors.

All together, these results indicate that our approach well captures quantitative and qualitative DNA methylation changes upon neuronal differentiation that might be underestimated or overlooked by an average methylation based approach.

DISCUSSION

Each cell is a unicum. Evidence has accumulated that even in morphologically homogeneous cell populations extensive differences can be highlighted at multiple molecular levels, and that these differences are relevant to biological processes (31).

Single-cell DNA methylation assays promise to be the standard technique to study DNA methylation heterogeneity in cell populations (30,62). However, single-cell DNA methylation technologies still generate very sparse data, in a limited number of cells per sample, and at high cost (28,32). Alternatively, cell-to-cell differences can be deduced by studying the methylation patterns of consecutive cytosines in sequenced reads from bulk experiments, assuming each read coming from a single DNA molecule (19,28,33–36). This approach can stand comparison with single-cell assays when the goal is to obtain a statistical/robust description of DNA methylation of a genomic region in a cell population (28).

Building on top of our experience on deep targeted bisulfite sequencing, in this study we propose MC profiling as a genome-wide approach to the study of DNA methylation heterogeneity. Given an epilocus holding 4 CpG sites, we defined its MC profile as the ensemble of the relative abundances of molecules sharing an equal number of methylated cytosines (Methylation Classes, MCs). Such an approach, while incorporating information on the overall average methylation of a region, directly informs on the different methylation levels, and their abundance, observed in a pool of molecules. This information is usually not, or poorly, taken into account by other approaches, which directly quantify the degree of cellular heterogeneity through the analysis of individual arrangements of methylated cytosines in single DNA molecules (epialleles).

A previous study showed that the methylation level of individual molecules can be used to adjust mean methylation indices for cell heterogeneity, thus improving prediction of gene expression levels compared to the overall average methylation (26). This method quantifies DNA methylation at promoter regions as the ratio of reads holding ≥ 1 methylated cytosines to the total number of reads mapped to the promoter. MC profiling is in line with this logic, but further enlarges the information on DNA methylation heterogeneity by considering molecules with different methylation levels as separate entities.

A conceptually similar approach to MC profiling has been proposed in (37–39). In these studies, DNA methylation is expressed as the probability mass function (PMF) of methylation levels that could be observed in a pool of molecules, which resemble the concept of our MCs. This approach, specifically designed to deal with the low coverage of WGBS experiments, has provided novel insights on DNA methylation heterogeneity and its disposition across the genome, its evolution upon differentia-

tion, aging and cancer, and its relationship with the genetic background (37–39). The biggest difference between this approach and MC profiling is that while the PMF is predicted from a mathematical model applied to DNA methylation data, the frequencies of MCs are empirically estimated from experimental data, thus avoiding time consuming model fitting and releasing the distribution of methylation from a-priori parametrization of DNA methylation dynamics.

To quantify the dissimilarity between MC profiles, we adopted the Jensen-Shannon distance (42). This dissimilarity measure has been applied in bioinformatics and epigenetics (37–39,63–65).

To set the parameters of our approach, we synthesized low coverage 4 CpG datasets from an in-house database of high-coverage amplicon bisulfite sequencing data (27,44–46). In this context, we provided a systematic quantification of the impact of coverage on the accuracy of MC profiles, and estimated the expected error associated with MC profiles at a coverage of 50 reads. In our opinion, these results could serve as guidelines to orient qualitative analysis of DNA methylation in low coverage settings.

Here, similarly to previous studies (40,41), we adopted a classification procedure, assigning each MC profile to the most similar among 5 reference profiles. This classification scheme provided us an interpretable representation of each MC profile. Furthermore, it provided us with a qualitative property to be compared across epiloci. Finally, being this a fixed scheme, we could apply it and directly compare the results on different conditions and species.

We demonstrated that MC profiles were stable among different samples and neighboring epiloci. Previous studies illustrated that DNA methylomes exhibit high inter-individual stability, especially in CG dense regions (66,67). Concordant epigenetic marks, including DNA methylation, across genomic blocks have been also described (23,37,68). Altogether, our results are in line with previously described patterns of regional and inter-individual stability of DNA methylation, and suggest that MC profiles capture controlled DNA methylation dynamics rather than stochastic fluctuations of methylation levels.

In this paper, we applied MC profiling to gain insights on methylation heterogeneity in various biological contexts.

Firstly, we profiled regions either present in single copies in the genome of individual cells or carrying heterozygous polymorphisms that enabled distinguishing the two alleles. We found that cell-to-cell differences were the strongest contributor to the molecular heterogeneity incorporated in MC profiles. Allelic differences contributed only in about 6% of analyzed regions, in agreement with the fraction of allele specific methylation described in most studies (38,49,69).

Secondly, we tested the capability of MC profiling to inspect known examples of mono-allelic regulation, i.e. genomic imprinting and X-inactivation, in which DNA methylation is notably involved.

When we analyzed the MC profiles of epiloci located in proximity of known genomic imprinted regions, we found that bimodal MC profiles were overrepresented. This was expected, considering the known opposite methylation pattern of the two parental alleles at imprinted regions (70). For an epilocus located upstream of the *Zdbf2* gene, holding a

polymorphic site, we were able to clearly show opposite MC profiles on the two alleles.

Loss of imprinting (LoI) has been described in several tumors (59,60). As a proof of concept, we showed that MC profiling can capture LoI in a leukemic sample, suggesting that this approach could be adopted in this field. We illustrated an example of MC profile alteration in the promoter of *GNAS*, which LoI has been described in diverse types of cancer (60). Consistently with previous study, we found that the MC profile was altered toward gain of DNA methylation in the tumor sample (60). In addition, MC profiling suggested that this gain was not homogeneously accomplished in the whole cell population.

We then analyzed the MC profiles of epiloci located on the X chromosome in female samples. First, we compared MC profiles of epiloci flanking genes with reported differential inactivation status. Consistent with previous findings (71–73), escapee epiloci showed homogeneous DNA methylation on both X copies, being unimodally fully methylated or unmethylated. Subject epiloci, on the contrary, were enriched for more heterogeneous MC profiles (D3 and D5), compatible with different DNA methylation status of the two alleles (71–73).

To further inspect the DNA methylation status of the inactive X, we selected the X epiloci with a fully unmethylated profile in male samples, and examined the corresponding MC profiles in female samples to infer the profile of the inactive X. We showed a prevalence of intermediately methylated classes on the inactive X, accompanied by high cellular heterogeneity. Incomplete DNA methylation of the inactive X was described in (74) at single CpG level, thus marking a difference between the X inactivation and the genomic imprinting processes that was well reflected in our analysis. It is worth noting that the prevalence of intermediately methylated MCs that we found with our approach also suggested a difference between the methylated status on the inactive X and at autosomal epiloci, suggestive of peculiar mechanisms intervening in DNA methylation establishment and regulation on the inactive X.

The methylation status of the inactive X appeared also to be highly heterogeneous among different cells. We speculate that this cellular epipolymorphism could almost in part find its reflection in differences of X inactivation status between equivalent cells described in single-cell RNA-seq studies (75,76).

Finally, we applied MC profiling to the analysis of DNA methylation changes in different conditions. In particular, we examined profiles' changes upon differentiation, when epigenetic remodeling is expected to occur. We adopted the Jensen-Shannon distance to capture epiloci with significant differences in MC profiles between neural precursors and differentiated neurons. Being JSD a symmetric distance measure, it did not return the information on whether MC profiles changes correspond to gain or loss of DNA methylation. Thus, we examined the pattern transitions to gain insights on how profiles' changes were occurring. Combining the analysis of JSD and pattern transitions provided us a comprehensive picture of DNA methylation differences among conditions: in fact, we could distinguish profiles changes associated with unvaried patterns (and thus, with stable reciprocal proportion of DNA molecules with differ-

ent methylation levels) from profiles changes accompanied with pattern transitions (which indicate a redistribution of the proportions of molecules with different methylation levels).

As expected, we found that MC profiles changes captured by JSD correlated with average DNA methylation gain or loss at most epiloci. However, we described MC profile changes at almost constant DNA methylation for >5000 epiloci. Qualitative DNA methylation changes occurring with little to no changes in overall average methylation were also described in (37–39), indicating that such an approach can be even more informative than average methylation based approach in the analysis of dynamic systems.

Interestingly, we found that MC profile changes were enriched at CpG islands, which were described to be spared from most epigenetic changes in the original study (50). The association that we found with changes of chromatin marks, as well as the concordance of MC profiles changes upon differentiation in two different neuronal subtypes, pointed to exclude random variations occurring at these epiloci. Instead, considering that most epiloci exhibited stable prototype classes in precursors and differentiated neurons, it is possible that MC profiling has captured changes in cellular heterogeneity that were overlooked by the average methylation-based approach.

Applying MC profiling to RRBS data can give insights on cellular epigenetic heterogeneity from plenty of already available datasets in public repositories. However, it strongly limits the analysis to CpG islands and immediately proximate regions (15,16). This limit is further exacerbated when selecting target regions harboring four CpGs (the epiloci of this study) shared among multiple samples. The required coverage of 50 reads strongly limits the applicability of the proposed approach outside Whole Genome Bisulfite Sequencing (WGBS) data. However, more unbiased enrichment assays have been developed which combine high throughput sequencing with selection of target regions through PCR or capture-based trapping that are natively less biased toward CG dense regions and could fit the coverage requirements of MC profiling (77,78).

Despite these limitations, we here showed that MC profiling could effectively capture cellular differences and changes also in CG dense regions, which are usually reported to be resistant to DNA methylation in normal conditions (3,79,80). We indeed believe that applying MC profiling to these experiments could further extend our observations outside CG dense regions.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

DATA AVAILABILITY

The datasets analyzed in this study are publicly available in GEO with accessions GSE66121, GSE130735, GSE53714, GSE72700. The R code utilized in this study has been deposited in Zenodo (<https://doi.org/10.5281/zenodo.7414513>).

FUNDING

No external funding

Conflict of interest statement. None declared.

REFERENCES

- Kim, M. and Costello, J. (2017) DNA methylation: an epigenetic mark of cellular memory. *Exp. Mol. Med.*, **49**, e322.
- Moore, L.D., Le, T. and Fan, G. (2013) DNA methylation and its basic function. *Neuropsychopharmacology*, **38**, 23–38.
- Jones, P.A. (2012) Functions of DNA methylation: islands, start sites, gene bodies and beyond. *Nat. Rev. Genet.*, **13**, 484–492.
- Greenberg, M.V.C. and Bourc'his, D. (2019) The diverse roles of DNA methylation in mammalian development and disease. *Nat. Rev. Mol. Cell Biol.*, **20**, 590–607.
- Smith, Z.D. and Meissner, A. (2013) DNA methylation: roles in mammalian development. *Nat. Rev. Genet.*, **14**, 204–220.
- Reik, W. and Walter, J. (2001) Genomic imprinting: parental influence on the genome. *Nat. Rev. Genet.*, **2**, 21–32.
- Suelves, M., Carrió, E., Núñez-Álvarez, Y. and Peinado, M.A. (2016) DNA methylation dynamics in cellular commitment and differentiation. *Brief. Funct. Genomics*, **15**, 443–453.
- Jin, Z. and Liu, Y. (2018) DNA methylation in human diseases. *Genes Dis.*, **5**, 1–8.
- Ehrlich, M. (2019) DNA hypermethylation in disease: mechanisms and clinical relevance. *Epigenetics*, **14**, 1141–1163.
- Robertson, K.D. (2005) DNA methylation and human disease. *Nat. Rev. Genet.*, **6**, 597–610.
- Nishiyama, A. and Nakanishi, M. (2021) Navigating the DNA methylation landscape of cancer. *Trends Genet.*, **37**, 1012–1027.
- Yong, W.-S., Hsu, F.-M. and Chen, P.-Y. (2016) Profiling genome-wide DNA methylation. *Epigenetics Chromatin*, **9**, 26.
- Masser, D.R., Berg, A.S. and Freeman, W.M. (2013) Focused, high accuracy 5-methylcytosine quantitation with base resolution by benchtop next-generation sequencing. *Epigenetics Chromatin*, **6**, 33.
- Varley, K.E. and Mitra, R.D. (2010) Bisulfite patch PCR enables multiplexed sequencing of promoter methylation across cancer samples. *Genome Res.*, **20**, 1279–1287.
- Gu, H., Smith, Z.D., Bock, C., Boyle, P., Gnirke, A. and Meissner, A. (2011) Preparation of reduced representation bisulfite sequencing libraries for genome-scale DNA methylation profiling. *Nat. Protoc.*, **6**, 468–481.
- Bock, C., Tomazou, E.M., Brinkman, A.B., Müller, F., Simmer, F., Gu, H., Jäger, N., Gnirke, A., Stunnenberg, H.G. and Meissner, A. (2010) Quantitative comparison of genome-wide DNA methylation mapping technologies. *Nat. Biotechnol.*, **28**, 1106–1114.
- Bock, C. (2012) Analysing and interpreting DNA methylation data. *Nat. Rev. Genet.*, **13**, 705–719.
- Elliott, G., Hong, C., Xing, X., Zhou, X., Li, D., Coarfa, C., Bell, R.J.A., Maire, C.L., Ligon, K.L., Sigaroudinia, M. *et al.* (2015) Intermediate DNA methylation is a conserved signature of genome regulation. *Nat. Commun.*, **6**, 6363.
- Scherer, M., Nebel, A., Franke, A., Walter, J., Lengauer, T., Bock, C., Müller, F. and List, M. (2020) Quantitative comparison of within-sample heterogeneity scores for DNA methylation data. *Nucleic Acids Res.*, **48**, e46.
- Irizarry, R.A., Ladd-Acosta, C., Wen, B., Wu, Z., Montano, C., Onyango, P., Cui, H., Gabo, K., Rongione, M., Webster, M. *et al.* (2009) The human colon cancer methylome shows similar hypo- and hypermethylation at conserved tissue-specific CpG island shores. *Nat. Genet.*, **41**, 178–186.
- Jaenisch, R. and Bird, A. (2003) Epigenetic regulation of gene expression: how the genome integrates intrinsic and environmental signals. *Nat. Genet.*, **33**, 245–254.
- Haerter, J.O., Lövkvist, C., Dodd, I.B. and Sneppen, K. (2014) Collaboration between CpG sites is needed for stable somatic inheritance of DNA methylation states. *Nucleic Acids Res.*, **42**, 2235–2244.
- Zhang, L., Xie, W.J., Liu, S., Meng, L., Gu, C. and Gao, Y.Q. (2017) DNA methylation landscape reflects the spatial organization of chromatin in different cells. *Biophys. J.*, **113**, 1395–1404.
- Lokk, K., Modhukur, V., Rajashekar, B., Märten, K., Mägi, R., Kolde, R., Koltšina, M., Nilsson, T.K., Vilo, J., Salumets, A. *et al.* (2014) DNA methylome profiling of human tissues identifies global and tissue-specific methylation patterns. *Genome Biol.*, **15**, 3248.
- Doi, A., Park, I.-H., Wen, B., Murakami, P., Aryee, M.J., Irizarry, R., Herb, B., Ladd-Acosta, C., Rho, J., Loewer, S. *et al.* (2009) Differential methylation of tissue- and cancer-specific CpG island shores distinguishes human induced pluripotent stem cells, embryonic stem cells and fibroblasts. *Nat. Genet.*, **41**, 1350–1353.
- Hansen, K.D., Timp, W., Bravo, H.C., Sabunciyan, S., Langmead, B., McDonald, O.G., Wen, B., Wu, H., Liu, Y., Diep, D. *et al.* (2011) Increased methylation variation in epigenetic domains across cancer types. *Nat. Genet.*, **43**, 768–775.
- Florio, E., Keller, S., Coretti, L., Affinito, O., Scala, G., Errico, F., Fico, A., Boscia, F., Sisalli, M.J., Reccia, M.G. *et al.* (2017) Tracking the evolution of epialleles during neural differentiation and brain development: d-Aspartate oxidase as a model gene. *Epigenetics*, **12**, 41–54.
- Huan, Q., Zhang, Y., Wu, S. and Qian, W. (2018) HeteroMeth: a database of Cell-to-cell heterogeneity in DNA methylation. *Genomics Proteomics Bioinformatics*, **16**, 234–243.
- Karemaker, I.D. and Vermeulen, M. (2018) Single-Cell DNA methylation profiling: technologies and biological applications. *Trends Biotechnol.*, **36**, 952–965.
- Angermueller, C., Clark, S.J., Lee, H.J., Macaulay, I.C., Teng, M.J., Hu, T.X., Krueger, F., Smallwood, S.A., Ponting, C.P., Voet, T. *et al.* (2016) Parallel single-cell sequencing links transcriptional and epigenetic heterogeneity. *Nat. Methods*, **13**, 229–232.
- Carter, B. and Zhao, K. (2021) The epigenetic basis of cellular heterogeneity. *Nat. Rev. Genet.*, **22**, 235–250.
- Teschendorff, A.E., Zhu, T., Breeze, C.E. and Beck, S. (2020) EPISCORE: cell type deconvolution of bulk tissue DNA methylomes from single-cell RNA-Seq data. *Genome Biol.*, **21**, 221.
- Landan, G., Cohen, N.M., Mukamel, Z., Bar, A., Molchadsky, A., Brosh, R., Horn-Saban, S., Zalcenstein, D.A., Goldfinger, N., Zundelovich, A. *et al.* (2012) Epigenetic polymorphism and the stochastic formation of differentially methylated regions in normal and cancerous tissues. *Nat. Genet.*, **44**, 1207–1214.
- Li, S., Garrett-Bakelman, F., Perl, A.E., Luger, S.M., Zhang, C., To, B.L., Lewis, I.D., Brown, A.L., D'Andrea, R.J., Ross, M.E. *et al.* (2014) Dynamic evolution of clonal epialleles revealed by methclone. *Genome Biol.*, **15**, 472.
- Xu, J., Shi, J., Cui, X., Cui, Y., Li, J.J., Goel, A., Chen, X., Issa, J.-P., Su, J. and Li, W. (2021) Cellular heterogeneity-adjusted clonal methylation (CHALM) improves prediction of gene expression. *Nat. Commun.*, **12**, 400.
- Zhang, X. and Wang, X. (2022) McConcord: a new metric to quantitatively characterize DNA methylation heterogeneity across reads and CpG sites. *Bioinformatics*, **38**, i307–i315.
- Jenkinson, G., Pujadas, E., Goutsias, J. and Feinberg, A.P. (2017) Potential energy landscapes identify the information-theoretic nature of the epigenome. *Nat. Genet.*, **49**, 719–729.
- Abante, J., Fang, Y., Feinberg, A.P. and Goutsias, J. (2020) Detection of haplotype-dependent allele-specific DNA methylation in WGBS data. *Nat. Commun.*, **11**, 5238.
- Jenkinson, G., Abante, J., Feinberg, A.P. and Goutsias, J. (2018) An information-theoretic approach to the modeling and analysis of whole-genome bisulfite sequencing data. *BMC Bioinformatics*, **19**, 87.
- Affinito, O., Scala, G., Palumbo, D., Florio, E., Monticelli, A., Miele, G., Avvedimento, V.E., Usiello, A., Chiariotti, L. and Coccozza, S. (2016) Modeling DNA methylation by analyzing the individual configurations of single molecules. *Epigenetics*, **11**, 881–888.
- De Riso, G., Fiorillo, D.F.G., Fierro, A., Cuomo, M., Chiariotti, L., Miele, G. and Coccozza, S. (2020) Modeling DNA methylation profiles through a dynamic equilibrium between methylation and demethylation. *Biomolecules*, **10**, 1271.
- Lin, J. (1991) Divergence measures based on the Shannon entropy. *IEEE Trans. Inf. Theory*, **37**, 145–151.
- Cover, T.M. and Thomas, J.A. (2005) In: *Elements of Information Theory*. 1st edn. Wiley.
- Cuomo, M., Keller, S., Punzo, D., Nuzzo, T., Affinito, O., Coretti, L., Carella, M., de Rosa, V., Florio, E., Boscia, F. *et al.* (2019) Selective demethylation of two CpG sites causes postnatal activation of the *dao*

- gene and consequent removal of d-serine within the mouse cerebellum. *Clin. Epigenetics*, **11**, 149.
45. Affinito, O., Palumbo, D., Fierro, A., Cuomo, M., De Riso, G., Monticelli, A., Miele, G., Chiariotti, L. and Coccozza, S. (2020) Nucleotide distance influences co-methylation between nearby CpG sites. *Genomics*, **112**, 144–150.
 46. Cuomo, M., Borrelli, L., Della Monica, R., Coretti, L., De Riso, G., D'Angelo Lancellotti di Durazzo, L., Fioretti, A., Lembo, F., Dinan, T.G., Cryan, J.F. *et al.* (2021) DNA methylation profiles of Tph1A and BDNF in gut and brain of l. Rhamnosus-treated zebrafish. *Biomolecules*, **11**, 142.
 47. Dahlet, T., Argüeso Lleida, A., Al Adhami, H., Dumas, M., Bender, A., Ngondo, R.P., Tanguy, M., Vallet, J., Auclair, G., Bardet, A.F. *et al.* (2020) Genome-wide analysis in the mouse embryo reveals the importance of DNA methylation for transcription integrity. *Nat. Commun.*, **11**, 3153.
 48. Kushwaha, G., Dozmorov, M., Wren, J.D., Qiu, J., Shi, H. and Xu, D. (2016) Hypomethylation coordinates antagonistically with hypermethylation in cancer development: a case study of leukemia. *Hum. Genomics*, **10**, 18.
 49. Urozco, L.D., Rubbi, L., Martin, L.J., Fang, F., Hormozdiari, F., Che, N., Smith, A.D., Lusk, A.J. and Pellegrini, M. (2014) Intergenerational genomic DNA methylation patterns in mouse hybrid strains. *Genome Biol.*, **15**, R68.
 50. Sharma, A., Klein, S.S., Barboza, L., Lohdi, N. and Toth, M. (2016) Principles governing DNA methylation during neuronal lineage and subtype specification. *J. Neurosci.*, **36**, 1711–1722.
 51. Krueger, F. and Andrews, S.R. (2011) Bismark: a flexible aligner and methylation caller for bisulfite-seq applications. *Bioinformatics*, **27**, 1571–1572.
 52. Sarnataro, A., De Riso, G., Coccozza, S., Pezone, A., Majello, B., Amente, S. and Scala, G. (2022) A novel workflow for the qualitative analysis of DNA methylation data. *Comput. Struct. Biotechnol. J.*, **20**, 5925–5934.
 53. Scala, G., Affinito, O., Palumbo, D., Florio, E., Monticelli, A., Miele, G., Chiariotti, L. and Coccozza, S. (2016) ampliMethProfiler: a pipeline for the analysis of CpG methylation profiles of targeted deep bisulfite sequenced amplicons. *BMC Bioinformatics*, **17**, 484.
 54. Krueger, F. and Andrews, S.R. (2016) SNPsplit: allele-specific splitting of alignments between genomes with known SNP genotypes. *F1000Research*, **5**, 1479.
 55. van der Velde, A., Fan, K., Tsuji, J., Moore, J.E., Purcaro, M.J., Pratt, H.E. and Weng, Z. (2021) Annotation of chromatin states in 66 complete mouse epigenomes during development. *Commun. Biol.*, **4**, 239.
 56. McLean, C.Y., Bristor, D., Hiller, M., Clarke, S.L., Schaar, B.T., Lowe, C.B., Wenger, A.M. and Bejerano, G. (2010) GREAT improves functional interpretation of cis-regulatory regions. *Nat. Biotechnol.*, **28**, 495–501.
 57. Landau, D.A., Clement, K., Ziller, M.J., Boyle, P., Fan, J., Gu, H., Stevenson, K., Sougnez, C., Wang, L., Li, S. *et al.* (2014) Locally disordered methylation forms the basis of intratumor methylome variation in chronic lymphocytic leukemia. *Cancer Cell*, **26**, 813–825.
 58. Hiura, H., Sugawara, A., Ogawa, H., John, R.M., Miyauchi, N., Miyanari, Y., Horiike, T., Li, Y., Yaegashi, N., Sasaki, H. *et al.* (2010) A tripartite paternally methylated region within the gpr1-zdbf2 imprinted domain on mouse chromosome 1 identified by meDIP-on-chip. *Nucleic Acids Res.*, **38**, 4929–4945.
 59. Plass, C. and Soloway, P.D. (2002) DNA methylation, imprinting and cancer. *Eur. J. Hum. Genet.*, **10**, 6–16.
 60. Kim, J., Bretz, C.L. and Lee, S. (2015) Epigenetic instability of imprinted genes in human cancers. *Nucleic Acids Res.*, **43**, 10689–10699.
 61. Balaton, B.P., Cotton, A.M. and Brown, C.J. (2015) Derivation of consensus inactivation status for X-linked genes from genome-wide studies. *Biol. Sex Differ.*, **6**, 35.
 62. Hui, T., Cao, Q., Wegrzyn-Woltosz, J., O'Neill, K., Hammond, C.A., Knapp, D.J.H.F., Laks, E., Moksa, M., Aparicio, S., Eaves, C.J. *et al.* (2018) High-Resolution single-cell DNA methylation measurements reveal epigenetically distinct hematopoietic stem cell subpopulations. *Stem Cell Rep.*, **11**, 578–592.
 63. Guo, X. (2020) JS-MA: a jensen-shannon divergence based method for mapping genome-wide associations on multiple diseases. *Front. Genet.*, **11**, 507038.
 64. Itzkovitz, S., Hodis, E. and Segal, E. (2010) Overlapping codes within protein-coding sequences. *Genome Res.*, **20**, 1582–1589.
 65. Kartal, Ö., Schmid, M.W. and Grossniklaus, U. (2020) Cell type-specific genome scans of DNA methylation divergence indicate an important role for transposable elements. *Genome Biol.*, **21**, 172.
 66. Palumbo, D., Affinito, O., Monticelli, A. and Coccozza, S. (2018) DNA methylation variability among individuals is related to cpGs cluster density and evolutionary signatures. *BMC Genomics*, **19**, 229.
 67. Bock, C., Walter, J., Paulsen, M. and Lengauer, T. (2008) Inter-individual variation of DNA methylation and its implications for large-scale epigenome mapping. *Nucleic Acids Res.*, **36**, e55.
 68. Ernst, J. and Kellis, M. (2017) Chromatin-state discovery and genome annotation with ChromHMM. *Nat. Protoc.*, **12**, 2478–2492.
 69. Onuchic, V., Lurie, E., Carrero, I., Pawliczek, P., Patel, R.Y., Rozowsky, J., Galeev, T., Huang, Z., Altschuler, R.C., Zhang, Z. *et al.* (2018) Allele-specific epigenome maps reveal sequence-dependent stochastic switching at regulatory loci. *Science*, **361**, eaar3146.
 70. Edwards, C.A. and Ferguson-Smith, A.C. (2007) Mechanisms regulating imprinted genes in clusters. *Curr. Opin. Cell Biol.*, **19**, 281–289.
 71. NISC Comparative Sequencing Program, Duncan, C.G., Grimm, S.A., Morgan, D.L., Bushel, P.R., Bennett, B.D., Roberts, J.D., Tyson, F.L., Merrick, B.A. and Wade, P.A. (2018) Dosage compensation and DNA methylation landscape of the x chromosome in mouse liver. *Sci. Rep.*, **8**, 10138.
 72. Cotton, A.M., Price, E.M., Jones, M.J., Balaton, B.P., Kobor, M.S. and Brown, C.J. (2015) Landscape of DNA methylation on the x chromosome reflects CpG density, functional chromatin state and X-chromosome inactivation. *Hum. Mol. Genet.*, **24**, 1528–1539.
 73. Balaton, B.P. and Brown, C.J. (2021) Contribution of genetic and epigenetic changes to escape from X-chromosome inactivation. *Epigenetics Chromatin*, **14**, 30.
 74. Balaton, B.P., Fornes, O., Wasserman, W.W. and Brown, C.J. (2021) Cross-species examination of X-chromosome inactivation highlights domains of escape from silencing. *Epigenetics Chromatin*, **14**, 12.
 75. Keniry, A. and Blewitt, M.E. (2018) Studying x chromosome inactivation in the single-cell genomic era. *Biochem. Soc. Trans.*, **46**, 577–586.
 76. Garieri, M., Stamoulis, G., Blanc, X., Falconnet, E., Ribaux, P., Borel, C., Santoni, F. and Antonarakis, S.E. (2018) Extensive cellular heterogeneity of x inactivation revealed by single-cell allele-specific expression in human fibroblasts. *Proc. Natl. Acad. Sci. U.S.A.*, **115**, 13015–13020.
 77. Kacmarczyk, T.J., Fall, M.P., Zhang, X., Xin, Y., Li, Y., Alonso, A. and Betel, D. (2018) “Same difference”: comprehensive evaluation of four DNA methylation measurement platforms. *Epigenetics Chromatin*, **11**, 21.
 78. Klobučar, T., Kreibich, E., Krueger, F., Arez, M., Pólvara-Brandão, D., von Meyenn, F., da Rocha, S.T. and Eckersley-Maslin, M. (2020) IMPLICON: an ultra-deep sequencing method to uncover DNA methylation at imprinted regions. *Nucleic Acids Res.*, **48**, e92.
 79. Edgar, R., Tan, P.P.C., Portales-Casamar, E. and Pavlidis, P. (2014) Meta-analysis of human methylomes reveals stably methylated sequences surrounding CpG islands associated with high gene expression. *Epigenetics Chromatin*, **7**, 28.
 80. Mohn, F., Weber, M., Rebhan, M., Roloff, T.C., Richter, J., Stadler, M.B., Bibel, M. and Schübeler, D. (2008) Lineage-Specific polycomb targets and de novo DNA methylation define restriction and potential of neuronal progenitors. *Mol. Cell*, **30**, 755–766.