

1 **Simultaneous measurement of chromatin accessibility, DNA methylation, and**
2 **nucleosome phasing in single cells**

3 Sebastian Pott¹

4 **1** University of Chicago, Department of Human Genetics, Chicago, IL, United States

5

6 Correspondence:

7 Sebastian Pott

8 University of Chicago

9 Department of Human Genetics

10 920 E. 58th Street, CLSC 317

11 Chicago, IL 60637

12 spott@uchicago.edu

13

14

15

16

17

18

19

1 **Introductory paragraph**

2 Gaining insights into the regulatory mechanisms that underlie the pervasive transcriptional variation
3 observed between individual cells^{1,2} necessitates the development of methods that measure
4 chromatin organization in single cells. *Nucleosome Occupancy* and *Methylome*-sequencing
5 (NOMe-seq) employs a GpC methyltransferase to detect accessible chromatin and has been used to
6 map nucleosome positioning and DNA methylation genome-wide in bulk samples^{3,4}. Here I provide
7 proof-of-principle that NOMe-seq can be adapted to measure chromatin accessibility and
8 endogenous DNA methylation in single cells (scNOMe-seq). scNOMe-seq recovered characteristic
9 accessibility and DNA methylation patterns at DNase hypersensitive sites (DHSs) and enabled
10 direct estimation of the fraction of accessible DHSs within an individual cell. In addition, scNOMe-
11 seq provided high resolution of chromatin accessibility within individual loci which was exploited
12 to detect footprints of CTCF binding and to estimate the average nucleosome phasing distances in
13 single cells. This approach could be applied to characterize the chromatin organization of single
14 cells in heterogeneous mixtures of cells, for example to samples of primary cancer cells.

15

16 **Main**

17 A number of methods that map chromatin organization in populations of cells previously have been
18 adapted for single cells, including ATAC-seq^{5,6}, DNase-seq⁷, methylome sequencing^{8,9}, and CHIP-
19 seq¹⁰. Interpretation of these data in single cells is complicated because the resulting signal is near
20 binary and extremely sparse^{5,7,11-13}. *Nucleosome Occupancy* and *Methylome*-sequencing (NOMe-
21 seq)³ employs the GpC methyltransferase (MTase) from *M.CviPI* to probe chromatin
22 accessibility^{3,4}. The GpC MTase methylates cytosines in GpC di-nucleotides in non-nucleosomal
23 DNA *in vitro*. Combined with high-throughput bisulfite sequencing this approach has been used to
24 characterize nucleosome positioning and endogenous methylation in human cell lines^{3,14} and in

1 selected promoters of single yeast cells¹⁵. NOMe-seq data have several unique features that are
2 advantageous in light of the challenges associated with single cell measurements (**Fig. 1a**). First,
3 NOMe-seq simultaneously measures chromatin accessibility (through GpC methylation) and
4 endogenous CpG DNA methylation. Chromatin accessibility indicates whether a putative regulatory
5 region might be utilized in a given cell¹⁶, while endogenous DNA methylation in regulatory regions
6 has been connected to a variety of regulatory processes often associated with repression¹⁷. The
7 ability to combine complementary assays within single cells is essential for a comprehensive
8 genomic characterization of individual cells since each cell represents a unique biological sample
9 which is almost inevitably destroyed in the process of the measurement. Second, each sequenced
10 read might contain several GpCs which independently report the accessibility status along the
11 length of that read. NOMe-seq therefore captures additional information compared to purely count-
12 based methods, such as ATAC-seq and DNase-seq, which increases the confidence associated with
13 the measurements and allows detection of footprints of individual transcription factor (TF) binding
14 events in single cells. Third, the DNA is recovered and sequenced independently of its methylation
15 status, which is a pre-requisite to distinguish between true negatives (i.e. closed chromatin) and
16 false negatives (i.e. loss of DNA) when assessing accessibility at specified locations in single cells.
17 NOMe-seq can therefore measure the fraction of accessible regions among a set of covered, pre-
18 defined genomic locations.

19 To adapt the NOMe-seq protocol^{3,18} to single cells (scNOMe-seq), individual nuclei were isolated
20 using fluorescence-activated cell sorting (FACS) and sorted into wells of a 96-well plate following
21 the incubation with the GpC MTase (**Fig. 1b and Supplemental Fig. 1**). DNA from isolated nuclei
22 was subjected to bisulfite conversion and sequencing libraries were prepared using a commercial kit
23 for amplification of low input bisulfite converted DNA (**Methods**). In this proof-of-concept study, I
24 used the well-characterized lymphoblast cell lines GM12878 and K562 to assess the feasibility and
25 performance of scNOMe-seq. The scNOMe-seq datasets in this study represent 19 individual
26 GM12878 cells and 12 individual K562 cells. The set of GM12878 cells included seven control

1 cells that were not incubated with GpC MTase (**Supplemental Fig. 2**). Each GpC MTase treated
2 library was sequenced to at least 16 M 100 bp reads, of which 37%- 64% aligned to the human
3 genome using the bisulfite aligner Bismark¹⁹ (**Supplemental Table 1**). Genome-wide, the number
4 of cytosines covered in GpC and CpG contexts averaged 6,679,864 (2.9%) of all GpCs and
5 1,291,180 (3.6%) of all CpGs per cell (**Supplemental Fig. 3 and Supplemental Table 1**).

6 To test whether the GpC methylation observed in GpC MTase treated samples (**Supplemental Fig.**
7 **4**) captured chromatin accessibility at specific genomic features, I focused on DNase hypersensitive
8 sites (DHSs) previously identified in GM12878 and K562 cell lines¹⁶. DHSs were associated with
9 strong enrichment of GpC methylation. This was observed both in data from pooled and individual
10 GM12878 cells (**Fig. 1 c, d and Supplemental Fig. 5**), and K562 cells (**Supplemental Fig. 6, 7**).

11 Conversely, endogenous DNA methylation decreased around the center of the DHSs in agreement
12 with previous reports^{20,21} (**Figure 1c and Supplemental Fig. 6**). These data show that scNOME-seq
13 detects chromatin accessibility at DHSs. In principle, the frequent occurrence of GpC di-nucleotides
14 renders the majority of DHSs detectable by NOME-seq (**Supplemental Fig. 8, 9**). On average in
15 10.6% (20388/191566) and 17.3% (33182/191598) of the DHSs with one or more GpC at least one
16 GpC was covered by a sequencing read, and in 5.2% (9083/174896) and 9.5% (16608/174828) of
17 the DHSs with four or more GpC at least four GpCs were covered in individual GM12878 cells and
18 K562 cells, respectively (**Fig. 1 e**). Chromatin accessibility signal can vary along the length of a
19 given DHSs due to binding of transcription factors²² and the specific position of a GpC within a
20 DHSs will thus affect its chance of being methylated. To account for this variability and to obtain
21 more robust estimates of GpC methylation only DHSs with at least 4 covered GpC were used for
22 the subsequent analyses and referred to as ‘covered DHSs’. Average GpC methylation of covered
23 DHSs in single cells was strongly associated with the observed DNaseI accessibility at these sites in
24 bulk populations (**Fig. 1f and Supplemental Fig. 10**). The opposite trend was observed for
25 endogenous CpG methylation which was lowest for DHSs with the highest DNaseI accessibility
26 (**Fig. 1g and Supplemental Fig. 10**). At the level of individual sites the distribution of GpC

1 methylation suggested that around 50% of the covered DHS showed no or low accessibility (i.e.
2 less than 25% GpC methylation) in individual cells (**Supplemental Fig. 11**). To estimate the
3 proportion of covered DHSs that were concurrently accessible in a single cells I applied a fixed
4 threshold of 40% GpC methylation above which sites were considered accessible (**Methods**). At
5 this GpC methylation threshold 32%-44% and 26%-37% of all covered DHSs were determined
6 accessible in single GM12878 and K562 cells, respectively. As expected these result depended to
7 some degree on the GpC methylation and the numbers of GpCs required to include a DHSs in the
8 analysis (**Supplemental Fig. 12**). However, even under the most lenient conditions less than 50%
9 of DHSs were accessible in most individual cells. Grouping the DHSs based on DNaseI
10 accessibility measured in bulk samples revealed that the degree of DNaseI accessibility closely
11 related to the frequency DHSs accessibility in single cells (**Fig. 1h**). This analysis leveraged the
12 NOME-seq-specific property that the DNA sequence is recovered independently of its accessibility
13 status. It provided direct evidence for the notion that the degree of DNaseI accessibility observed in
14 DNase-seq of bulk samples reflects the frequency with which a particular region is accessible in
15 individual cells.

16 A potentially powerful application for single cell genomic approaches is the label-free classification
17 of single cells from heterogeneous mixtures of cells solely based on the measured feature^{5,11,23}. Of
18 note, using a union set of DHSs from both cell types was sufficient to classify individual GM12878
19 and K562 cells into their respective cell types based on GpC methylation (**Fig. 1 i**). While this
20 assessment might have been influenced in part by the separate processing of both cell types, both
21 cell types showed preferential enrichment of GpC methylation at their respective DHSs compared
22 to DHSs identified in the other cell type (**Supplemental Fig. 13**). Thus, this approach should be
23 extendable to scNOME-seq data from samples containing mixtures of cell types and endogenous
24 CpG methylation could be included in such analyses to provide additional information⁹.

25 To examine in detail whether scNOME-seq captures features of chromatin accessibility that are
26 specifically associated with transcription factor binding I analyzed scNOME-seq data at

1 transcription factor binding sites (TFBS). The average GpC methylation around CTCF ChIP-seq
2 peaks¹⁶ in single cells recapitulated the accessibility previously observed in NOMe-seq bulk
3 samples³: Accessibility increased strongly towards the CTCF binding sites while the location of the
4 CTCF motif at the center of the region showed low accessibility suggesting that CTCF binding
5 protected from GpC MTase activity and thus creating a footprint of CTCF binding, both when
6 averaged across data from all single cells (**Fig. 2a and Supplemental Fig. 14**) and in individual
7 cells (**Fig. 2b and Supplemental Fig. 15**). In contrast, endogenous CpG methylation was generally
8 depleted around the center of CTCF binding sites (**Fig. 2a and Supplemental Fig. 14**). Similar
9 accessibility profiles, albeit less pronounced compared to CTCF were observed for additional
10 transcription factors (**Supplemental Fig. 16**). These analyses provided evidence, that in aggregate,
11 scNOMe-seq detected CTCF DNA binding events from single cells. In addition to aggregated
12 binding sites, scNOMe-seq data should also provide information about transcription factor binding
13 at individual loci. To test whether scNOMe-seq data detected CTCF footprints at individual motifs,
14 GpC methylation at motifs within CTCF ChIP-seq peaks that contained at least one GpC was
15 compared to the GpC methylation level in the regions flanking each motif (**Fig. 2c**). On average,
16 two-thirds of CTCF motif instances within these accessible regions showed no GpC methylation, or
17 lower GpC methylation than the flanking regions suggesting that scNOMe-seq detected footprints
18 caused by binding of CTCF (**Fig. 2d and f**). Of note, motifs associated with a footprint had
19 significantly higher scores than motifs without a footprint suggesting that the motif score is a strong
20 determinant of CTCF binding within these accessible regions (**Fig. 2e, g and Supplemental Fig.**
21 **17**). The CTCF footprint could be observed at individual loci and comparing GpC methylation from
22 multiple cell and suggests that scNOMe-seq also detected cell-to-cell variation in the footprint
23 (**Figure 2h and Supplemental Fig. 18**). Footprint measurements such as this should be generally
24 feasible for TFs whose motifs contain at least one GpC di-nucleotide and could be used to infer the
25 activity of a wide range of transcription factors in single cells.

1 The pattern of GpC methylation adjacent to CTCF sites suggested that scNOME-seq also detected
2 the well-positioned nucleosomes flanking these regions (**Fig. 2a**)³. This observation was confirmed
3 by the oscillatory distribution of the average GpC and CpG methylation around locations of well-
4 positioned nucleosomes identified from MNase-seq data¹⁶ (**Fig. 3a and Supplemental Fig. 14**).
5 While nucleosome core particles are invariably associated with DNA fragments of 147 bp,
6 nucleosomes are separated by linker DNA of varying lengths, resulting in different packaging
7 densities between cell types²⁴ and between genomic regions within a cell^{24,25}. To determine whether
8 scNOME-seq data can be used to measure the average linker length, average distances between
9 nucleosome midpoints in single cells (phasing distances) were estimated by correlating the
10 methylation status between pairs of cytosines in GpC di-nucleotides at offset distances from 3 bp to
11 400 bp (**Fig. 3c, d and Supplemental Fig. 19, 20**). The estimated phases fell between 187 bp and
12 196 bp (mean=196.7 bp) in GM12878 cells, and between 188 bp and 200 bp (mean =194.2 bp) in
13 K562 cells (**Fig. 3e**). These estimates are in general agreement with phase estimates derived from
14 MNase-seq data in human cells²⁴. In addition, estimated phasing distances varied within individual
15 cells depending on the chromatin context, similar to observation from bulk MNase-seq data²⁴ (**Fig.**
16 **3f**). These proof-of-principle experiments have been performed using commercial kits for bisulfite
17 conversion and library amplification, additional optimization or alternative amplification
18 approaches⁸ are likely to increase the yield substantially. Ultimately, it should be possible to
19 integrate the GpC MTase treatment into microfluidic workflows and combine this method with
20 scRNA-seq, similar to recently published methods that combine scRNA-seq and methylome-
21 sequencing²⁶. scNOME-seq will be particularly useful for studies that aim to simultaneously
22 measure chromatin accessibility and DNA methylation, or that aim to measure activity of
23 transcription factors in individual cells based on footprints in single cells. This approach could be
24 applied to characterize the chromatin organization of single cells in heterogeneous mixtures of cells,
25 for example to samples of primary cancer cells. scNOME-seq could be applied to characterize the

1 chromatin organization of single cells in heterogeneous mixtures of cells, for example in samples of
2 primary cancer cells.

3

4 **Methods**

5 **Cell culture, nuclei isolation, and GpC methylase treatment**

6 GM12878 and K562 cells were obtained from Coriell and ATCC, respectively. GM12878 were
7 grown in RPMI medium 1640 (Gibco), supplemented with 2mM L-Glutamine (Gibco), and
8 Penicilin and Streptavidin (Pen Strep, Gibco), and 15% fetal bovine serum (FBS, Gibco). K562
9 were grown in RPMI medium 1640 of the same composition but with 10% FBS. Cells were grown
10 at 37 C and in 5% CO₂. NOME-Seq procedure was performed based on protocols for CpG
11 methyltransferase SSSI described in¹⁸ and the GpC methyltransferase from *M.CviPI*³, with some
12 modification. Between 2x10⁶ and 5x10⁶ cells were harvested by centrifuging the cell suspension
13 for 5 min at 500x g. Cells were washed once with 1x PBS, re-suspended in 1 ml lysis buffer (10mM
14 Tris-HCl pH 7.4, 10mM NaCl, 3mM MgCl₂) and incubated for 10 min on ice. IGEPAL CA-630
15 (Sigma) was added to a final concentration of 0.025% and the cell suspension was transferred to a 2
16 ml Dounce homogenizer. Nuclei were released by 15 strokes with the pestle. Success of lysis was
17 confirmed by inspection under a light microscope. Nuclei were collected by centrifuging the cell
18 suspension for 5 min at 800x g at 4C and washed twice with cold lysis buffer without detergent.
19 One million nuclei were resuspended in reaction buffer to yield a suspension with a final
20 concentration of 1x GpC MTase buffer (NEB), 0.32 mM S-Adenosylmethionine (SAM) (NEB), and
21 50 ul of GpC methyltransferase (4U/ul) from *M.CviPI* (NEB). The suspension was carefully mixed
22 before incubating for 8 min at 37 C after which another 25 ul of enzyme and 0.7 ul of 32 mM SAM
23 were added for an additional 8 min incubation at 37C. To avoid disruption of nuclei incubation was
24 stopped by adding 750 ul of 1x PBS and collecting the nuclei at 800 xg. Supernatant was removed

1 and nuclei were re-suspended in 500ul 1x PBS containing Hoechst 33342 DNA dye (NucBlue Live
2 reagent, Hoechst). Nuclei were kept on ice until sorting.

3 **Nuclei isolation using Fluorescence activated cell sorting (FACS), lysis, and DNA bisulfite** 4 **conversion**

5 Nuclei were sorted at the Flow Cytometry core at the University of Chicago on a BD FACSAria or
6 BD FACSAria Fusio equipped with a 96-well-plate holder. To obtain individual and intact nuclei
7 gates were set on forward and side scatter to exclude aggregates and debris. DAPI/PacBlue channel
8 or Violet 450/500 channel were used to excite the Hoechst 33342 DNA dye and to gate on cells
9 with DNA content corresponding to cells in G1 phase of the cell cycle in order to maintain similar
10 DNA content per cell and to remove potential heterogeneity attributable to cell cycle. Cells were
11 sorted into individual wells pre-filled with 19 ul of 1x M-Digestion buffer (EZ DNA Methylation
12 Direct Kit, Zymo Research) containing 1 mg/ml Proteinase K. Following collection, the plates were
13 briefly spun to collect droplets that might form during handling. Nuclei were lysed by incubating
14 the samples at 50 C for 20 min in a PCR cycler. DNA was subjected to bisulfite conversion by
15 adding 130 ul of freshly prepared CT Conversion reagent (EZ DNA Methylation Direct Kit, Zymo)
16 to the lysed nuclei. Conversion was performed by denaturing the DNA at 98 C for 8 min followed
17 by 3.5 hrs incubation at 65 C. DNA isolation was performed using the EZ DNA Methylation Direct
18 Kit (Zymo Research) following the manufacturer's instruction with the modification that the DNA
19 was eluted in only 8 ul of elution buffer.

20 **Library preparation and sequencing**

21 Libraries were prepared using the Pico Methyl-seq Library prep Kit (Zymo Research) following the
22 manufacturer's instruction for low input samples. Specifically, the random primers were diluted 1:2
23 before the initial pre-amplification step and the first amplification was extended to a total of 10
24 amplification cycles. Libraries were amplified with barcoded primers allowing for multiplexing.
25 The sequences can be found in **Supplemental Table 2**, primers were ordered from IDT. The
26 purification of amplified libraries was performed using Agencourt AMPureXP beads (Beckmann

1 Coulter), using a 1:1 ratio of beads and libraries. Concentration and size distribution of the final
2 libraries was assessed on an Bioanalyzer (Agilent). Libraries with average fragment size above 150
3 bp were pooled and sequenced. Libraries were sequenced on Illumina HiSeq 2500 in rapid mode
4 (K562 cells) and HiSeq4000 (GM12878 cells).

5 **Read processing and alignment**

6 Sequences were obtained using 100 bp paired-end mode. For processing and alignment each read
7 from a read pair was treated independently as this slightly improved the mapping efficiency. Before
8 alignment, read sequences in fastq format were assessed for quality using fastqc
9 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Reads were trimmed to remove low
10 quality bases and 6 bp were clipped from the 5 prime end of each read to avoid mismatches
11 introduced by amplification. In the case of GM12878 cells 6 bp were clipped from either end of the
12 read. Only reads that remained longer than 20 bp were kept for further analyses. These processing
13 steps were performed using trim_galore version 0.4.0
14 (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) with the following settings:
15 *trim_galore --quality 30 --phred33 --illumina --stringency 1 -e 0.1 --clip_R1 6 --gzip --length 20 --*
16 *output_dir outdir Sample.fastq.gz*. The trimmed fastq files were aligned using the bisulfite aligner
17 bismarck version 0.15.0²⁷ which calls bowtie2²⁸ internally. Reads were aligned to the human
18 genome (genome assembly hg38). Reads were aligned in single read mode using default settings.
19 The amplification protocol used to generate the scNOME-seq libraries yielded non-directional
20 libraries and alignment was performed with the option `--non_directional` (*bismark --fastq --prefix*
21 *SamplePrefix --output_dir output_dir --non_directional --phred33-quals --score_min L,0,-0.2 --*
22 *bowtie2 genome_file trimmed.fastq.gz*). Some libraries contained small amounts of DNA from *C.*
23 *elegans* as spike-ins, however these were not used during the analysis. Duplicates were removed
24 using samtools version 0.1.19²⁹ on sorted output files from bismark (*samtools rmdup*
25 *SamplePrefix.sorted.bam SampleAligned_rmdup.bam*).

26 **Extraction of GpC and CpG methylation status**

1 Coverage and methylation status of all cytosines was extracted using
2 `bismark_methylation_extractor`²⁷ (`bismark_methylation_extractor -s --ignore 6 --output outdir --`
3 `cytosine_report --CX --genome_folder path_to_genome_data SampleAligned_rmdup.bam`). The
4 resulting coverage files were used to extract the methylation status of cytosines specifically in GpC
5 and CpG di-nucleotides using the `coverage2cytosine` script which is part of Bismark²⁷. The resulting
6 coverage files contained cytosines in GCG context which are ambiguous given that they represent a
7 cytosine both in GpC and CpG di-nucleotides. Coordinates of these ambiguous positions were
8 identified using `oligoMatch`³⁰ and these positions were removed from the coverage files. The
9 number of unconverted cytosines (estimated based on apparent methylation rates in non-GpC and
10 non-CpG context) was low in all libraries (<1%). However, it was noted that unconverted cytosines
11 were not randomly distributed but associated with entirely unconverted reads. Regions covered by a
12 read with more than 3 unconverted cytosines in non-CpG and non-GpC context were removed from
13 further analysis as well. The genotype was not taken into account as its effect on calling the
14 methylation status incorrectly was deemed negligible for the analyses performed here.

15 **Analysis of GpC and CpG methylation at genomic features in single cells**

16 ScNOME-seq data were compared to a number of genomic features in GM12878 and K562 cells
17 collected by ENCODE²⁴ which were downloaded through the UCSC data repository³¹. These
18 datasets are listed in **Supplemental Table 3**. While the scNOME-seq data were aligned against
19 human genome assembly hg38, some of the datasets were only available on genome assembly hg19
20 and the coordinates of these datasets were lifted from hg19 to hg38 using `liftOver`³⁰ (default re-
21 mapping ratio 1). Nucleosome positions based on MNase-seq data in GM12878 were determined
22 with DANPOS version 2.2.2³² using default settings. Resulting intervals were lifted to hg38. After
23 removing summit locations with occupancy values above 300, the top 5% (713361) of nucleosome
24 positions based on their summit occupancy value were used.

25 GpC and CpG methylation density across intervals encompassing DNase hypersensitivity sites
26 (DHSs), transcription factor binding sites (TFBS), and well positioned nucleosomes was calculated

1 across the 2 kb regions centered on the middle of these regions using the scoreMatrixBin function in
2 the genomation package³³ in R³⁴. Data were aggregated in 5 bp bins for each region and across all
3 regions covered in a single cell. The average methylation level in pre-defined intervals (DHSs,
4 TFBS) was determined by computing the average GpC or CpG methylation for each interval
5 together with the number of GpC/CpGs covered in this interval using the map function from the
6 bedtools³⁵ suite. If not specifically mentioned otherwise DHSs were considered ‘covered’ and used
7 in analyses when at least 4 GpCs occurring within the predefined interval were covered by
8 sequencing data in an individual cell. Covered DHSs that Because the frequency of CpG di-
9 nucleotides is significantly lower, only 2 CpGs were required in order for a DHSs to be considered
10 covered for analyses focused on endogenous DNA methylation. To estimate the number of
11 cytosines within a given DHSs that could be covered only cytosines on the forward strand were
12 counted. While each GpC dinucleotide can be measured on both strands and would therefore yield a
13 count of two cytosines the data are sparse and each location will get at most a single read. This
14 approach should therefore give a more conservative estimate of the possible GpC coverage. For
15 analyses that used the scores of the peak regions, the peak scores reported the datasets from bulk
16 samples were used¹⁶.

17 For analyses that were centered on transcription factor binding motifs the PWMs were obtained
18 from the JASPAR database (2014)³⁶ for the TFs CTCF (MA0139), EBF1 (MA0154), and
19 PU.1(MA0080). Genome-wide scanning for locations of sequence matches to the PWMs was
20 performed using matchPWM in the Biotstring package³⁷ in R with a threshold of 75% based on the
21 human genome assembly hg38.

22 All plots were prepared using ggplot2³⁸, with the exception of heatmaps displaying the average
23 methylation density around genomic features in individual cells which were prepared using
24 heatmap.2 in gplots³⁹.

25 **Comparison of chromatin accessibility between cells**

1 Similarity in accessible chromatin between cells was calculated based on Jaccard similarity. Jaccard
2 similarity index (eq. 1) was calculated between pairs of samples by first obtaining the intersection
3 of DHSs covered in both samples of a pair with more than 4 GpCs. Each features was annotated as
4 open or closed, depending on the methylation status ($\geq 40\%$ methylation) and only pairs in which
5 at least one of the members was open.

$$6 \quad \text{jac}(A, B) = \frac{(A \cap B)}{(A \cup B)} \quad (1)$$

7 The similarity between samples from GM12878 and K562 cells was calculated based on the union
8 set of DHSs from both cell lines. The similarity indexes of all pairwise comparisons were used to
9 compute the distances between each cell. The resulting clustered data were displayed as a heat map.

10 **CTCF footprints in single cells**

11 CTCF footprints were measured by comparing the GpC methylation level in each motif to the
12 methylation level in the 50bp flanking regions immediately upstream and downstream of the motif.
13 Overlapping motifs were merged into a single interval before determining the coordinates for
14 flanking regions. To ensure sufficient GpC coverage for each interval the resulting three interval
15 was required to have at least one covered GpC and 4 GpCs covered in total among the three
16 intervals. This analysis only included regions that were accessible based on the methylation status
17 of the flanking regions (at least 50%). A CTCF footprint 'score' was determined by simply
18 subtracting the average GpC methylation of the flanking regions from the GpC methylation of the
19 motif.

20 scNOME-seq data were displayed in the UCSC genome browser³⁰ by converting the GpC
21 methylation coverage file into a bed file and using the methylation value as score. To facilitated the
22 visualization of the data in the context of previous Encode data the methylation files were lifted to
23 hg19. The tracks shown together with scNOME-seq data are Open Chromatin by DNaseI HS from
24 ENCODE/OpenChrom (Duke University) for DNaseI hypersensitivity, Nucleosome Signal from

1 ENCODE/Stanford/BYU, and CTCF ChIP-seq signal from Broad Histone Modification by ChIP-
2 seq from ENCODE/Broad Institute. All data are from GM12878 cells.

3 **Estimation of nucleosome phasing**

4 Nucleosome phasing estimates were obtained by first calculating the correlation coefficients for the
5 methylation status of pairs of GpCs at different offset distances. These values were computed using
6 a custom python script. Essentially, pairs of sequenced cytosines in GpC di-nucleotides were
7 collected for each offset distance from 3bp to 400bp cytosine. At each offset distance the correlation
8 of the methylation status was calculated across all pairs. Correlation coefficients were plotted
9 against the offset distances revealing periodic changes in the correlation coefficient. The
10 smoothed data were used to estimate the phasing distances by obtaining the offset distance
11 corresponding to the local maximum found between 100 bp and 300 bp. To determine phase lengths
12 of nucleosomes in different chromatin contexts the GpC coverage files were filtered for positions
13 falling into categories defined by chromHMM^{16,40} before obtaining the correlation coefficients.

14 **Data access**

15 Raw data and methylation coverage files are available at GEO (<https://www.ncbi.nlm.nih.gov/geo/>)
16 under the accession number . Reviewers might use this link:
17 <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=glotcwqqjbqlvef&acc=GSE83882>

18

19 **Competing financial interest**

20 The author declares no competing financial interests

21

22

23 **Acknowledgements**

1 I like to thank Yoav Gilad for support, and Greg Crawford and colleagues in the Department of
2 Human Genetics for helpful suggestions and comments on the manuscript. Cell sorting was
3 performed by M. Olson and D. Leclerc at the Flow Cytometry core of the University of Chicago. I
4 am grateful to Jason Lieb for input and support at the beginning of this project.

5

6 **Bibliography**

- 7 1. Shalek, A. K. *et al.* Single-cell transcriptomics reveals bimodality in expression and splicing in
8 immune cells. *Nature* **498**, 236–240 (2013).
- 9 2. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-
10 cell RNA-seq. *Nature* **509**, 371–375 (2014).
- 11 3. Kelly, T. K. *et al.* Genome-wide mapping of nucleosome positioning and DNA methylation
12 within individual DNA molecules. *Genome Research* **22**, 2497–2506 (2012).
- 13 4. Kilgore, J. A., Hoose, S. A., Gustafson, T. L., Porter, W. & Kladde, M. P. Single-molecule and
14 population probing of chromatin structure using DNA methyltransferases. *Methods (San Diego,
15 Calif.)* **41**, 320–332 (2007).
- 16 5. Cusanovich, D. A. *et al.* Multiplex single cell profiling of chromatin accessibility by
17 combinatorial cellular indexing. *Science* **348**, 910–914 (2015).
- 18 6. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory
19 variation. *Nature* 1–15 (2015). doi:10.1038/nature14590
- 20 7. Jin, W. *et al.* Genome-wide detection of DNase I hypersensitive sites in single cells and FFPE
21 tissue samples. *Nature* 1–17 (2015). doi:10.1038/nature15740

- 1 8.Smallwood, S. A. *et al.* Single-cell genome-wide bisulfite sequencing for assessing epigenetic
2 heterogeneity. *Nature Methods* **11**, 817–820 (2014).
- 3 9.Farlik, M. *et al.* Single-Cell DNA Methylome Sequencing and Bioinformatic Inference of
4 Epigenomic Cell-State Dynamics. *CellReports* **10**, 1386–1397 (2015).
- 5 10.Rotem, A. *et al.* Single-cell ChIP-seq reveals cell subpopulations defined by chromatin state.
6 *Nature Biotechnology* **33**, 1–11 (2015).
- 7 11.Buenrostro, J. D., Wu, B., Chang, H. Y. & Greenleaf, W. J. ATAC-seq: A Method for Assaying
8 Chromatin Accessibility Genome-Wide. *Current protocols in molecular biology / edited by*
9 *Frederick M. Ausubel ... [et al.]* **109**, 21.29.1 (2015).
- 10 12.Maurano, M. T. & Stamatoyannopoulos, J. A. Taking Stock of Regulatory Variation. *Cell*
11 *Systems* **1**, 18–21 (2015).
- 12 13.Bernstein, B. E. *et al.* Charting a dynamic DNA methylation landscape of the human genome.
13 *Nature* 1–5 (2013). doi:10.1038/nature12433
- 14 14.Taberlay, P. C., Statham, A. L., Kelly, T. K., Clark, S. J. & Jones, P. A. Reconfiguration of
15 nucleosome-depleted regions at distal regulatory elements accompanies DNA methylation of
16 enhancers and insulators in cancer. *Genome Research* **24**, 1421–1432 (2014).
- 17 15.Small, E. C., Xi, L., Wang, J.-P., Widom, J. & Licht, J. D. Single-cell nucleosome mapping
18 reveals the molecular basis of gene expression heterogeneity. *Proceedings of the National Academy*
19 *of Sciences of the United States of America* **111**, NaN–NaN (2014).
- 20 16.ENCODE Project Consortium & The ENCODE Project Consortium. An integrated
21 encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- 22 17.Schübeler, D. Function and information content of DNA methylation. *Nature* **517**, 321–326
23 (2015).

- 1 18.Miranda, T. B., Kelly, T. K., Bouazoune, K. & Jones, P. A. Methylation-Sensitive Single-
2 Molecule Analysis of Chromatin Structure. 1–16 (2001). doi:10.1002/0471142727.mb2117s89
- 3 19.Krueger, F., Kreck, B., Franke, A. & Andrews, S. R. DNA methylome analysis using short
4 bisulfite sequencing data. *Nature Methods* **9**, 145–151 (2012).
- 5 20.Stadler, M. B. *et al.* DNA-binding factors shape the mouse methylome at distal regulatory
6 regions. *Nature* 1–7 (2011). doi:10.1038/nature10716
- 7 21.Ziller, M. J. *et al.* Dissecting neural differentiation regulatory networks through epigenetic
8 footprinting. *Nature* **518**, 355–359 (2015).
- 9 22.Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor
10 footprints. *Nature* **489**, 83–90 (2012).
- 11 23.Jaitin, D. A. *et al.* Massively parallel single-cell RNA-seq for marker-free decomposition of
12 tissues into cell types. *Science* **343**, 776–779 (2014).
- 13 24.Valouev, A. *et al.* Determinants of nucleosome organization in primary human cells. *Nature*
14 **474**, 516–520 (2011).
- 15 25.Schones, D. E. *et al.* Dynamic Regulation of Nucleosome Positioning in the Human Genome.
16 *Cell* **132**, 887–898 (2008).
- 17 26.Angermueller, C. *et al.* Parallel single-cell sequencing links transcriptional and epigenetic
18 heterogeneity. *Nature Methods* 1–6 (2016). doi:10.1038/nmeth.3728
- 19 27.Krueger, F. & Andrews, S. R. Bismark: a flexible aligner and methylation caller for Bisulfite-
20 Seq applications. *Bioinformatics* **27**, 1571–1572 (2011).
- 21 28.Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**,
22 357–359 (2012).

- 1 29.Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–
2 2079 (2009).
- 3 30.Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Research* **12**, 996–1006
4 (2002).
- 5 31.Karolchik, D. *et al.* The UCSC Genome Browser database: 2014 update. *Nucleic Acids*
6 *Research* **42**, NaN–NaN (2014).
- 7 32.Chen, K. *et al.* DANPOS: dynamic analysis of nucleosome position and occupancy by
8 sequencing. *Genome Research* **23**, 341–351 (2013).
- 9 33.Akalin, A., Franke, V., Vlahoviček, K., Mason, C. E. & Schübeler, D. Genomation: a toolkit to
10 summarize, annotate and visualize genomic intervals. *Bioinformatics* **31**, 1127–1129 (2015).
- 11 34.R Core Team. R: A language and environment for statistical computing. (2015). at
12 <<https://www.R-project.org/>>
- 13 35.Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
14 features. *Bioinformatics* **26**, 841–842 (2010).
- 15 36.Tan, G. JASPAR2014: Data package for JASPAR. at <<http://jaspar.genereg.net/>>
- 16 37.Pages, H., Aboyoun, P., Gentleman, R. C. & DebRoy, S. Biostrings: String objects representing
17 biological sequences, and matching algorithms. (2016).
- 18 38.Wickham, H. ggplot2. 213 (2009). doi:10.1007/978-0-387-98141-3
- 19 39.Warnes, G. R., Bolker, B., Bonebakker, L. & Gentleman, R. gplots: Various R programming
20 tools for plotting data. (2016). at <<https://CRAN.R-project.org/package=gplots>>
- 21 40.Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types.
22 *Nature* **473**, 43–49 (2011).

1 **Figure Legends**

2 **Figure 1. scNOME-seq detected DNase hypersensitive sites in single cells.** a) Schematic of GpC
3 methyltransferase-based mapping of chromatin accessibility and simultaneous detection of
4 endogenous DNA methylation. b) Schematic of scNOME-seq procedure introduced in this study. c)
5 Average GpC methylation level (blue) and CpG methylation level (orange) at DNase
6 Hypersensitive sites (DHSs) in GM12878 cells. Regions are centered on the middle of DNase-seq
7 peak locations. Shown is the average methylation across a 2 kb window of 12 GM12878 cells. d)
8 Heatmap displaying the average GpC methylation level across the same regions as in c). Each row
9 corresponds to an individual GM12878 cell. Cells were grouped by similarity. e) Proportion of
10 DHSs covered by scNOME-seq data in each cell. The proportion displayed corresponds to the
11 fraction of DHSs covered by at least 1 or 4 GpCs in a given cell. Only DHSs with at least 1 or 4
12 GpCs, respectively, within their primary sequence were taken in consideration. Error bars represent
13 the standard deviation. f) Average GpC methylation and g) endogenous CpG methylation at DHSs
14 split into 10 groups based on associated DNase-seq peak scores from lowest to highest scores.
15 Average methylation per DHSs was based on the data from 12 GM12878 cells. h) Fraction of
16 accessible sites in individual GM12878 and K562 cells, respectively. DHSs were grouped by
17 DNase-seq peak scores. DHSs was considered accessible if the average methylation for that locus
18 was above 40%. Only DHSs with at least 4 covered GpCs were included. i) Heatmap shows
19 similarity scores (pair-wise jaccard distances) between all GM12878 and K562 cells. Comparison
20 was based on the union set of DHSs from GM12878 and K562 cells. Cells were grouped based on
21 unsupervised hierarchical clustering.

22

23

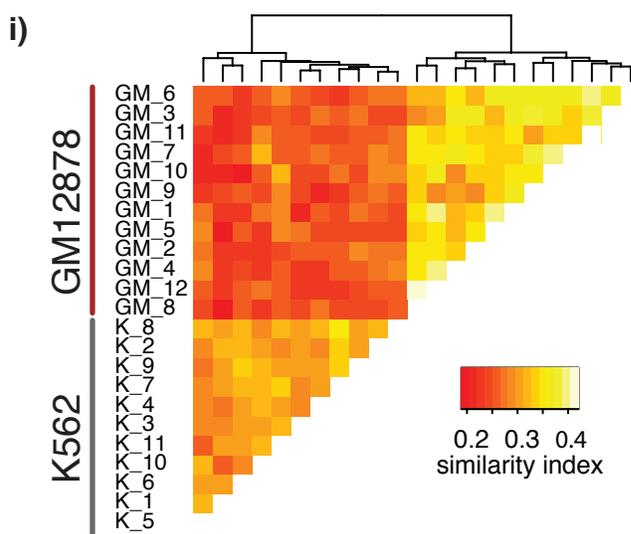
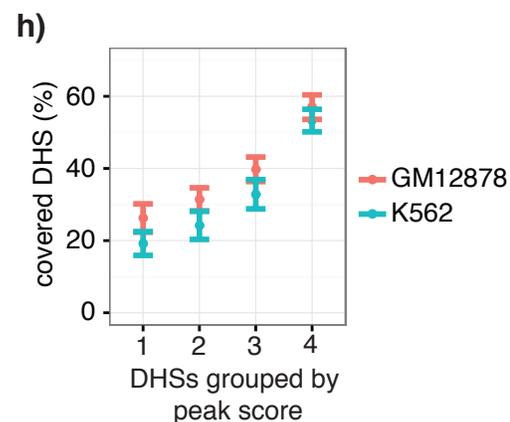
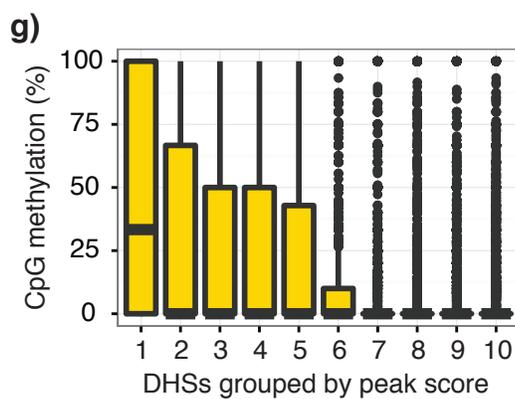
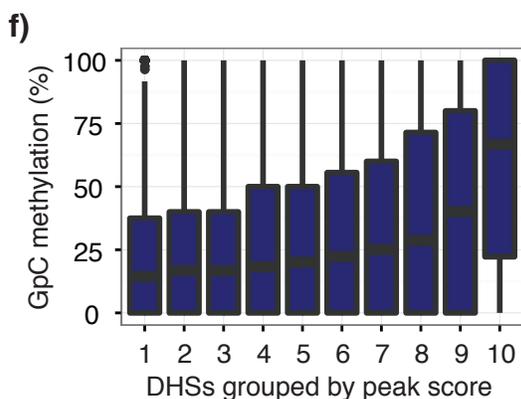
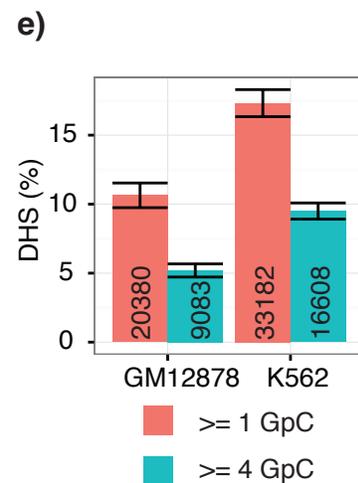
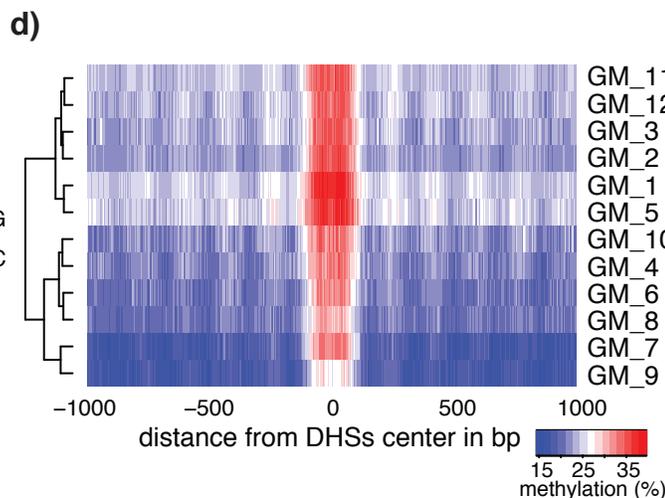
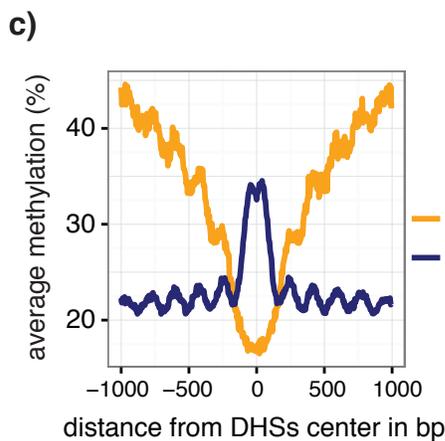
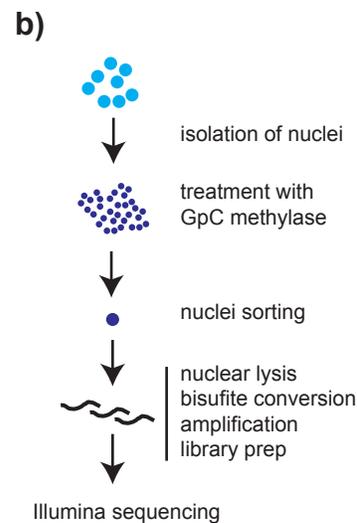
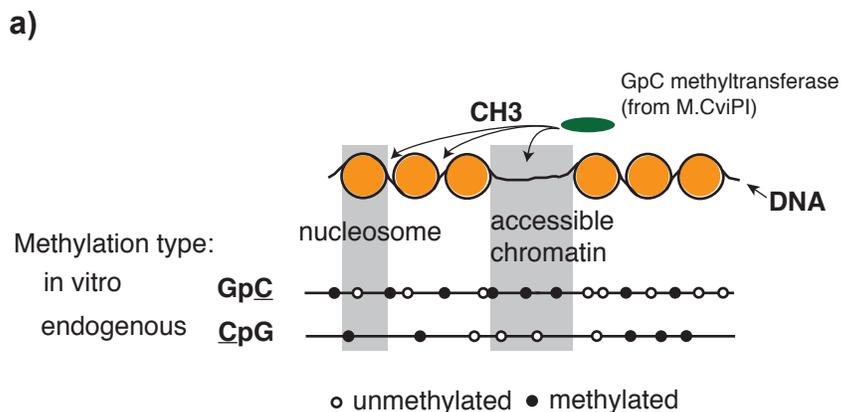
24

1 **Figure 2. scNOME-seq detected characteristic accessibility patterns at CTCF transcription**
2 **factor binding sites and measured CTCF footprints at individual loci** a) Average GpC
3 methylation level (blue) and CpG methylation level (orange) at CTCF binding sites in GM12878
4 cells. Regions are centered on motif locations. Shown is the average methylation across a 2 kb
5 window of the pool of 12 GM12878 cells. b) Heatmap displaying the average GpC methylation
6 across CTCF binding sites. Each row corresponds to an individual GM12878 cell and rows are
7 grouped by similarity. c) Schematic outline the measurement of CTCF footprints in accessible
8 regions. M denotes CTCF binding motifs within CTCF ChIP-seq regions and U and D indicate 50
9 bp upstream and downstream flanking regions. footprint score was determined by subtracting the
10 average GpC methylation in the flanking regions from the GpC methylation at the motif. d)
11 Heatmap displays GpC methylation in accessible regions found in a representative GM12878 cell
12 (GM_1). Each row represents a single CTCF motif instance within a CTCF ChIP-seq region.
13 Average methylation values for the motif and the 50 bp upstream and downstream regions are
14 shown separately. Regions are sorted based on the footprint score. Displayed are only regions that
15 had sufficient GpC coverage and that were considered accessible based on the methylation status of
16 the flanking regions. e) Heatmap reporting the CTCF motif scores for the motif regions in d).
17 Regions are sorted in the same order as in d). f) Average number of accessible regions at CTCF
18 motifs and the average number of those with a detectable footprint per individual GM12878 cell.
19 Error bars reflect standard deviation. g) Average CTCF motif scores in regions with and without
20 CTCF footprint for all 12 GM12878 cells. Each line connects the two data points from an individual
21 cell h) Combined display of scNOME-seq data from this study and DNase hypersensitivity data,
22 nucleosome occupancy, and CTCF ChIP-seq data from ENCODE. Upper panel shows a ~10 kb
23 region containing a CTCF binding sites. DNaseI hypersensitivity data and nucleosome density show
24 characteristic distribution around CTCF binding sites in GM12878 cells. Lower panel shows the
25 GpC methylation data of 5 individual cells that had sequencing coverage in this region, 4 of the
26 cells provide GpC data covering the CTCF motif located in the region. scNOME-seq data tracks

1 show methylation status of individual GpCs. Each row corresponds to data from a single cell. These
2 data indicate that binding of CTCF is detected in all 4 cells. Data are displayed as tracks in the
3 UCSC genome browser (<http://genome.ucsc.edu>).

4

5 **Figure 3. Nucleosome phasing in single cells.** a) Average GpC methylation level and b) CpG
6 methylation level at well-positioned nucleosomes in GM12878 cells. Regions are centered on
7 midpoints of top 5% of positioned nucleosomes. Shown is the average methylation across a 2 kb
8 window of the pool of 12 GM12878 cells. c), d) Correlation coefficients for the comparison in
9 methylation status between GpCs separated by different offset distances for GM12878 (c) and K562
10 (d) cells. Each line represents a single cell. Data are smoothed for better visualization. e)
11 Distribution of estimated phase lengths for GM12878 and K562 cells. f) Nucleosome phasing in
12 GM12878 in genomic regions associated with different chromatin states defined by chromHMM
13 (ENCODE). Boxplot represents the distribution of estimated phase lengths from all 12 GM12878
14 cells and overlaid points indicate values of each individual cells.



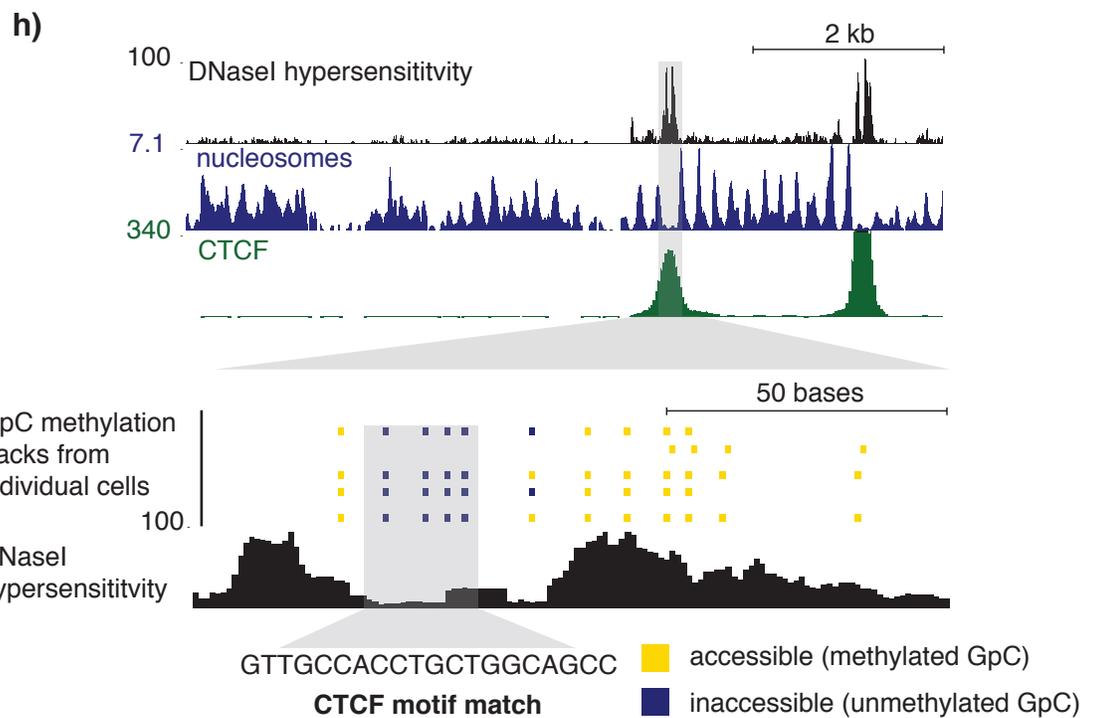
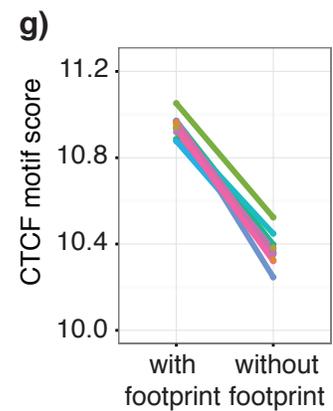
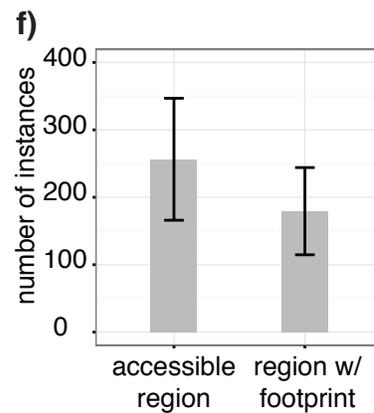
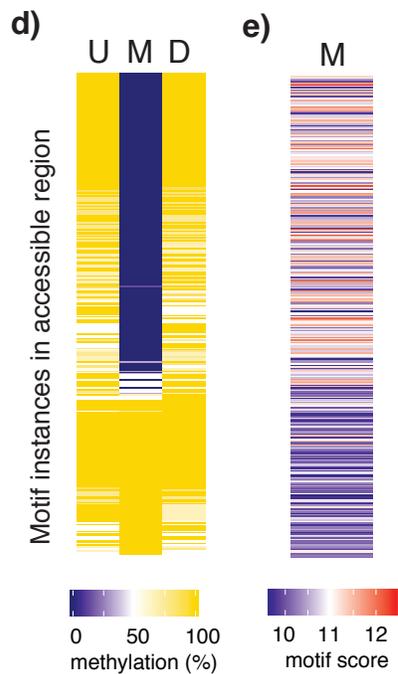
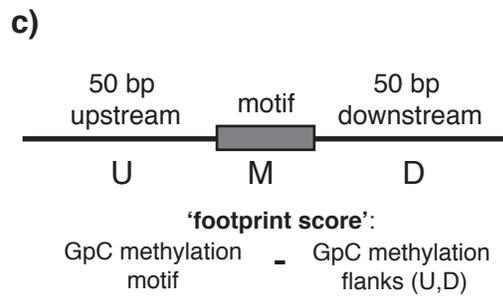
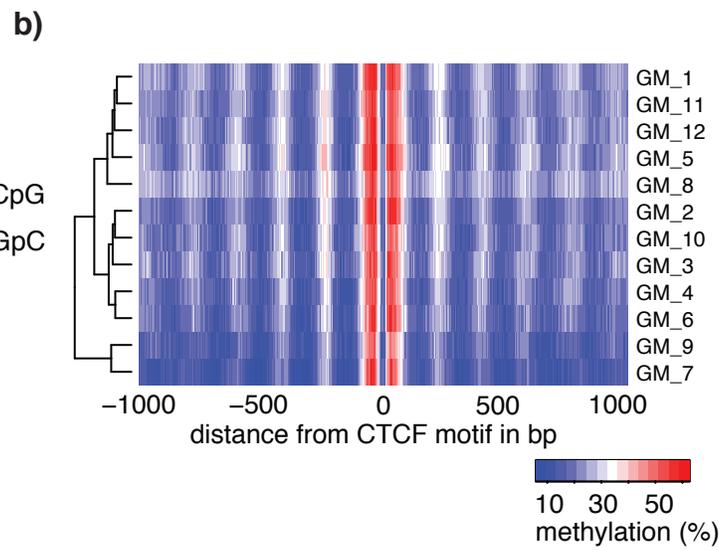
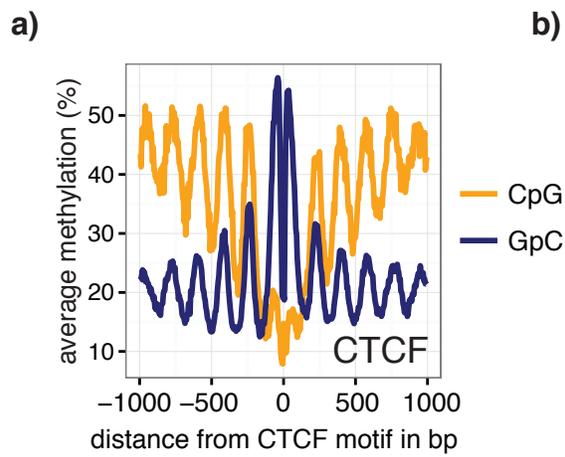


Figure 3