1    **Simultaneous measurement of chromatin accessibility, DNA methylation, and**

2    **nucleosome phasing in single cells**

3    Sebastian Pott

4    University of Chicago, Department of Human Genetics, Chicago, Il, United States

5

6    Correspondence:

7    Department of Human Genetics

8    University of Chicago

9    920 East 58$^{th}$ Street

10   Chicago, IL, 60637

11   Fax: (773) 834-0505

12   E-mail: spott@uchicago.edu

13

14

15

16

17

18

19   Running title: Chromatin organization in single cells

# 1    Abstract

2    Gaining insights into the regulatory mechanisms that underlie the transcriptional variation observed

3    between individual cells necessitates the development of methods that measure chromatin

4    organization in single cells. Here I adapted *N*ucleosome *O*ccupancy and *Me*thylome-sequencing

5    (NOMe-seq) to measure chromatin accessibility and endogenous DNA methylation in single cells

6    (scNOMe-seq). scNOMe-seq recovered characteristic accessibility and DNA methylation patterns

7    at DNase hypersensitive sites (DHSs). An advantage of scNOMe-seq is that sequencing reads are

8    sampled independently of the accessibility measurement. scNOMe-seq therefore controlled for

9    fragment loss, which enabled direct estimation of the fraction of accessible DHSs within individual

10    cells. In addition, scNOMe-seq provided high resolution of chromatin accessibility within

11    individual loci which was exploited to detect footprints of CTCF binding events and to estimate the

12    average nucleosome phasing distances in single cells. scNOMe-seq is therefore well-suited to

13    characterize the chromatin organization of single cells in heterogeneous cellular mixtures.

14

15

16

17

18

19

20

21

22

2

# 1 Introduction

2 Extensive transcriptional variation between individual cells has been observed using single cell

3 RNA-seq. These data facilitate identification of functional subpopulations in seemingly

4 homogeneous cell populations (Shalek et al. 2014), or characterization of the cellular composition

5 of complex tissues (Jaitin et al. 2014; Treutlein et al. 2014; Macosko et al. 2015). To gain

6 mechanistic insights into regulatory features that underlie cellular heterogeneity it is essential to

7 measure chromatin organization in individual cells. A number of methods that map chromatin

8 organization in populations of cells have been adapted for single cells, including ATAC-seq

9 (Cusanovich et al. 2015; Buenrostro et al. 2015b), DNase-seq (Jin et al. 2015), methylome

10 sequencing (Smallwood et al. 2014; Farlik et al. 2015), and ChIP-seq (Rotem et al. 2015).

11 Interpretation of these data in single cells is complicated because of the near binary and extremely

12 sparse signal (Cusanovich et al. 2015; Buenrostro et al. 2015b; Maurano and Stamatoyannopoulos

13 2015). *N*ucleosome *O*ccupancy and *Me*thylome-sequencing (NOMe-seq) (Kelly et al. 2012)

14 employs the GpC methyltransferase (MTase) from *M.CviPI* to probe chromatin accessibility (Kelly

15 et al. 2012; Kilgore et al. 2007). The GpC MTase methylates cytosines in GpC dinucleotides in

16 non-nucleosomal DNA *in vitro*. Combined with high-throughput bisulfite sequencing this approach

17 has been used to characterize nucleosome positioning and endogenous methylation in human cell

18 lines (Kelly et al. 2012; Taberlay et al. 2014) and in selected promoters of single yeast cells (Small

19 et al. 2014). NOMe-seq data have several unique features that are advantageous considering the

20 challenges associated with single cell measurements **(Fig. 1 a)**. First, NOMe-seq simultaneously

21 measures chromatin accessibility (through GpC methylation) and endogenous CpG methylation.

22 Chromatin accessibility indicates whether a putative regulatory region might be utilized in a given

23 cell (ENCODE Project ConsortiumThe ENCODE Project Consortium 2012), while endogenous

24 DNA methylation in regulatory regions has been connected to a variety of regulatory processes

25 often associated with repression (Schübeler 2015). The ability to combine complementary assays

26 within single cells is essential for a comprehensive genomic characterization of individual cells

3

1    since each cell represents a unique biological sample which is almost inevitably destroyed in the

2    process of the measurement. Second, each sequenced read might contain several GpCs which

3    independently report the accessibility status along the length of that read. NOMe-seq therefore

4    captures additional information compared to purely count-based methods, such as ATAC-seq and

5    DNase-seq, which increases the confidence associated with the measurements and allows detection

6    of footprints of individual transcription factor (TF) binding events in single cells. Third, the DNA is

7    recovered and sequenced independently of its methylation status, which is a pre-requisite to

8    distinguish between true negatives (i.e. closed chromatin) and false negatives (i.e. loss of DNA)

9    when assessing accessibility at specified locations in single cells. This is especially important in

10   single cells where allelic drop-out is pervasive. In single cells, NOMe-seq can therefore measure the

11   fraction of accessible regions among a set of covered, pre-defined genomic locations. In this proof-

12   of-principle study, I showed that NOMe-seq, which previously had only been performed on bulk

13   samples (Kelly et al. 2012; Taberlay et al. 2014), can be performed on single cells. In addition to

14   endogenous methylation at CpG dinucleotides, single cell NOMe-seq (scNOMe-seq) measured

15   chromatin accessibility at DHSs and TF binding sites in individual cells, and detected footprints of

16   CTCF binding at individual loci. Finally, the average phasing distance between nucleosomes within

17   individual cells can also be estimated from scNOMe-seq data.

18

## Results

20   To adapt the NOMe-seq protocol (Kelly et al. 2012; Miranda et al. 2010) to single cells, individual

21   nuclei were first incubated with GpC MTase and then sorted into wells of a 96-well plate using

22   fluorescence-activated cell sorting (FACS) (**Fig. 1b and Figure 1 – figure supplement 1**). DNA

23   from isolated nuclei was subjected to bisulfite conversion and sequencing libraries were prepared

24   using a commercial kit for amplification of low amounts of bisulfite-converted DNA (**Methods**). To

25   assess the feasibility and performance of NOMe-seq in single cells, I used the well-characterized

4

1   cell lines GM12878 and K562. The scNOMe-seq datasets in this study represent 19 individual

2   GM12878 cells and 11 individual K562 cells. The set of GM12878 cells included seven control

3   cells that were not treated with GpC MTase (**Figure 1– figure supplement 2**). Each GpC MTase-

4   treated library was sequenced to at least 16 M individual reads (**Methods).** Reads were aligned to

5   the human genome using the aligner Bismark (Krueger et al. 2012) and, after removal of duplicate

6   reads, between 2.5M and 5M reads were retained per library (**Supplementary file 1**). On average

7   6,679,864 (2.9%) of all cytosines in GpCs and 1,291,180 (3.6%) of all cytosines in CpGs were

8   covered per cell (**Figure 2– figure supplement 1 and Supplementary file 1**).

9

10  **scNOMe-seq accurately detected accessible chromatin at DNaseI hypersensitive sites**

11  To test whether the GpC methylation observed in GpC MTase treated samples (**Figure 2– figure**

12  **supplement 1)** captured known chromatin accessibility patterns, I focused on DNaseI

13  hypersensitive sites (DHSs) that were previously identified in GM12878 and K562 cell lines

14  (ENCODE Project ConsortiumThe ENCODE Project Consortium 2012). DHSs were associated

15  with strong enrichment of GpC methylation, both in data from pooled and individual GM12878

16  (**Figure 2 a, b, Figure 2– figure supplement 2**) and K562 cells (**Figure 2– figure supplement 3,**

17  **4**). Conversely, endogenous CpG methylation decreased around the center of the DHSs in

18  agreement with previous reports (Stadler et al. 2011; Ziller et al. 2014) (**Figure 2 a and Figure 2–**

19  **figure supplement 3**). These data show that scNOMe-seq detected chromatin accessibility at

20  DHSs. To assess how many of the DHSs regions were covered in a single cell, I first filtered DHSs

21  that contained GpC dinucleotides within their primary sequence and thus could be theoretically

22  detected by NOMe-seq. The frequent occurrence of GpC di-nucleotides renders the majority (>

23  85%) of DHSs detectable by NOMe-seq (**Figure 2– figure supplement 5, 6**). Of the theoretically

24  detectable DHSs, 10.6% (20388/191566) and 17.3% (33182/191598) had 1 or more GpCs covered

25  and, using a more stringent criterion, 5.2% (9083/174896) and 9.5% (16608/174828) were covered

1    at 4 or more GpCs in individual GM12878 cells and K562 cells, respectively (**Fig. 2 c**). Chromatin

2    accessibility signal can vary along the length of a given DHSs due to binding of transcription

3    factors (Neph et al. 2012) and the specific position of a GpC within a DHS will thus affect its

4    chance of being methylated. To account for this variability and to obtain more robust estimates of

5    GpC methylation only DHSs with at least 4 covered GpC were used for the subsequent analyses and

6    referred to as 'covered DHSs'.

7    In single cells, the average GpC methylation at covered DHSs was strongly correlated with the

8    observed DNaseI accessibility at these sites in bulk populations (**Fig. 2 d**, **Figure 2 –figure**

9    **supplement 7, 8)**. The opposite trend was observed for endogenous CpG methylation which was

10   lowest for DHSs with the highest DNaseI accessibility (**Figure 2 –figure supplement 7).** The

11   correlation between GpC methylation and DNaseI accessibility was lower for scNOMe-seq data

12   compared to bulk NOMe-seq data in the same cell line **(Figure 2 –figure supplement 8)**. At the

13   level of individual sites the distribution of GpC methylation suggested that around 50% of the

14   covered DHS showed less than 25% GpC methylation in individual cells (**Figure 2 –figure**

15   **supplement 9)**. To estimate the proportion of covered DHSs that were concurrently accessible in a

16   single cell I applied a fixed threshold of 40% GpC methylation above which sites were considered

17   accessible (**Methods**). At this GpC methylation threshold 32-44% and 26-37% of all covered DHSs

18   were determined to be accessible in single GM12878 and K562 cells, respectively. As expected

19   these results depended to some degree on the cutoffs used for GpC methylation and the number of

20   required GpCs per DHS. However, even under the most lenient conditions less than 50% of DHSs

21   were accessible in individual cells (**Figure 2 –figure supplement 10**). Grouping the DHSs based on

22   DNaseI accessibility in bulk samples, confirmed that the degree of DNaseI accessibility related

23   closely to the frequency of DHS accessibility in single cells (**Fig. 2 e**). This analysis leveraged the

24   NOMe-seq-specific property that the DNA sequence is recovered independently of its accessibility

25   status. It provided direct evidence for the notion that the degree of DNaseI accessibility observed in

26   DNase-seq of bulk samples reflects the frequency with which a region is accessible in individual

6

1    cells. Consequently, chromatin accessibility between cells is less variable at regions with high

2    DNaseI accessibility in bulk samples (**Figure 2 –figure supplement 11**). Correspondingly,

3    correlation of GpC methylation between individual cells is stronger at DHS loci compared to

4    randomized locations (**Figure 2 –figure supplement 12)**.

5    **scNOMe-seq captured characteristic chromatin organization associated with transcription**

6    Chromatin accessibility and endogenous methylation show characteristic patterns at gene promoters

7    and within gene bodies (Schübeler 2015; ENCODE Project ConsortiumThe ENCODE Project

8    Consortium 2012). To test whether these features can be observed in scNOMe-seq data, I first

9    plotted the average GpC and CpG methylation around transcription start sites (TSS). The average

10   GpC methylation showed the expected increase of chromatin accessibility directly upstream of the

11   TSS (**Fig. 3 a, Figure 3 – figure supplement 1**). In contrast, and as expected, the endogenous CpG

12   methylation decreased towards the TSS (**Fig. 3 b**). To visualize the distribution of CpG methylation

13   throughout entire gene loci, I plotted the aggregated CpG methylation across regions containing the

14   entire gene body and 50 kb upstream and 50 kb downstream of each gene (**Fig. 3 c, Figure 3 –**

15   **figure supplement 1**). Endogenous methylation was specifically reduced at the narrow promoter

16   region and gradually increased throughout the gene body. Downstream of the transcription end site

17   (TES) the average level CpG methylation level fell back to the non-genic background level.

18   Endogenous CpG methylation is typically increased within highly expressed genes (Schübeler

19   2015). This trend was clearly apparent in the single cell data where gene body methylation was

20   highest in highly expressed genes (**Fig. 3 d, Figure 3 –figure supplement 1**). Correspondingly, in

21   promoter regions (-500bp to +150bp) chromatin accessibility (GpC methylation) increased with the

22   transcript level of the adjacent gene (**Fig. 3 e, Figure 3 –figure supplement 2**). In contrast to

23   chromatin accessibility, endogenous methylation was lowest in promoters of genes with high

24   transcript levels (**Fig. 3 f**). These data show that scNOMe-seq recapitulated known characteristics of

25   chromatin accessibility and endogenous methylation at gene promoters and within gene bodies.

1    **GpC methylation and endogenous CpG methylation data separated individual GM12878 and**

2    **K562 cells**

3    A potentially powerful application for single cell genomic approaches is the label-free classification

4    of single cells from heterogeneous mixtures of cells solely based on the measured feature

5    (Cusanovich et al. 2015; Buenrostro et al. 2015a; Jaitin et al. 2014; Macosko et al. 2015). Of note,

6    using a union set of DHSs from both cell types was sufficient to classify individual GM12878 and

7    K562 cells into their respective cell types based on GpC methylation (**Fig. 4 a, Figure 4 –figure**

8    **supplement 1)**. While this assessment might have been influenced in part by the separate

9    processing of the cell types, both cell types showed preferential enrichment of GpC methylation at

10    their respective DHSs compared to DHSs identified in the other cell type **(Fig. 4 b)**. Similar to GpC

11    methylation, endogenous CpG methylation at multiple sets of genomic features was sufficient to

12    separate the cells into the respective cell types (**Fig. 4 c, Figure 4 –figure supplement 1)**.

13    **Detection of footprints of CTCF binding at individual loci in single cells**

14    To examine in detail whether scNOMe-seq captures features of chromatin accessibility that are

15    specifically associated with transcription factor binding, I analyzed scNOMe-seq data at

16    transcription factor binding sites (TFBS). The average GpC methylation around CTCF ChIP-seq

17    peaks (ENCODE Project ConsortiumThe ENCODE Project Consortium 2012) in single cells

18    recapitulated the accessibility previously observed in NOMe-seq bulk samples (Kelly et al. 2012):

19    Accessibility increased strongly towards the CTCF binding sites while the location of the CTCF

20    motif at the center of the region showed low accessibility suggesting that CTCF binding protected

21    from GpC MTase activity and thus creating a footprint of a CTCF binding event, both when

22    averaged across data from all single cells (**Fig. 5 a and Figure 5 – figure supplement 1**) and in

23    individual cells (**Fig. 5 b and Figure 5 – figure supplement 2**). In contrast, endogenous CpG

24    methylation was generally depleted around the center of CTCF binding sites **(Fig. 5 a and Figure 5**

25    **– figure supplement 1)**. Similar accessibility profiles, albeit less pronounced compared to CTCF,

8

1    were observed for additional transcription factors, for example EBF1 and PU.1 (**Figure 5 – figure**

2    **supplement 3**). These analyses provided evidence that, in aggregate, scNOMe-seq detected

3    chromatin accessibility characteristic of CTCF binding in single cells. To test whether scNOME-seq

4    data detected CTCF footprints at individual motifs loci, GpC methylation at motifs within CTCF

5    ChIP-seq peaks was compared to the GpC methylation level in the regions flanking each motif **(Fig.**

6    **5 c)**. On average, two-thirds of CTCF motif instances within these accessible regions showed no

7    GpC methylation, suggesting that CTCF binding prevented the GpC MTase from methylating the

8    cytosines within the binding motif and thus creating a footprint **(Fig. 5 d and f)**. Of note, motifs

9    associated with a footprint had significantly higher scores than motifs without a footprint suggesting

10   that the motif score is a strong determinant of CTCF binding within these accessible regions **(Fig. 5**

11   **e, g and Figure 5 – figure supplement 4)**. Of note, the CTCF footprints could be observed at

12   individual loci within individual cells and were shared across cells **(Figure 5 h and Figure 5 –**

13   **figure supplement 5).**

14

15   **Estimating nucleosome phasing in single cells**

16   The pattern of GpC methylation adjacent to CTCF sites suggested that scNOMe-seq also detected

17   the well-positioned nucleosomes flanking these regions (**Fig. 5 a**) (Kelly et al. 2012). This

18   observation was confirmed by the oscillatory distribution of the average GpC and CpG methylation

19   around locations of well-positioned nucleosomes identified from MNase-seq data (ENCODE

20   Project ConsortiumThe ENCODE Project Consortium 2012) (**Fig. 6 a**). While nucleosome core

21   particles are invariably associated with DNA fragments of  147 bp, nucleosomes are separated by

22   linker DNA of varying lengths, resulting in different packaging densities between cell types and

23   between genomic regions within a cell (Valouev et al. 2011; Schones et al. 2008). To determine

24   whether scNOMe-seq data can be used to measure the average linker length, average distances

25   between nucleosome midpoints in single cells (phasing distances) were estimated by correlating the

9

1    methylation status between pairs of cytosines in GpC di-nucleotides at offset distances from 3 bp to

2    400 bp **(Fig. 6 c, d and Figure 6 – figure supplement 1, 2).** The estimated phases fell between 187

3    bp and 196 bp (mean = 196.7 bp) in GM12878 cells, and between 188 bp and 200 bp (mean = 194.2

4    bp ) in K562 cells (**Fig. 6 e**). These estimates are in general agreement with phase estimates derived

5    from MNase-seq data in human cells (Valouev et al. 2011). In addition, estimated phasing distances

6    varied within individual cells depending on the chromatin context, similar to observation from bulk

7    MNase-seq data (Valouev et al. 2011) (**Fig. 6 f**).

8    **Discussion**

9    In this study, I demonstrated that scNOMe-seq simultaneously measures chromatin accessibility by

10   GpC methylation as well as endogenous CpG and DNA methylation in single cells. scNOMe-seq

11   detected chromatin accessibility at DHSs and TFBS and, in aggregate, these data recapitulated

12   NOMe-seq data obtained from bulk cells (Kelly et al. 2012). scNOMe-seq data also detected

13   footprints of CTCF binding, and was used to estimate nucleosome phasing distances.

14   Similar to other single cell genomic methods, scNOMe-seq relies on annotations obtained from bulk

15   measurements ((Cusanovich et al. 2015; Buenrostro et al. 2015b; Smallwood et al. 2014; Farlik et

16   al. 2015). A limitation of single cell genomic methods is their sparse coverage which leads to high

17   allelic drop-out. For methods in which the signal is based on counting the sequenced fragments,

18   such as ATAC-seq and DNase-seq, this poses a challenge since true negatives at a specific location

19   cannot be distinguished from false negatives that are a consequence of read loss. Compared to these

20   methods, scNOMe-seq has the unique advantage, that reads are recovered independently of the

21   signal and allelic drop-out events therefore can be distinguished from closed or inaccessible

22   chromatin configurations. The frequency of accessible sites in the population of DHSs can be

23   estimated. Using this approach only about 30-50 % of DHSs detected in the population were found

24   accessible in a single cell, depending on the thresholds chosen to call a site accessible. While this

25   assessment would have been possible using bulk NOMe-seq data, scNOMe-seq offers important

10

1   possibilities for future applications. For example, to compare accessibility across multiple loci

2   within a single cell and the use of heterogeneous cellular mixtures as input material.

3   As expected, the chance of a covered DHS to being open or closed is not equally distributed across

4   all DHSs from the population. Instead, DHSs with strong DNaseI accessibility showed a higher

5   frequency of accessibility in single cells compared to those sites with low DNaseI accessibility in

6   the population (**Fig. 2 e**) suggesting that the peak height is indeed directly related to the frequency

7   with which a site is accessible in individual cells. In agreement with this observation a large

8   proportion of variability observed between cells was attributable to DHSs with low DNaseI

9   accessibility in bulk samples (**Figure2 − figure supplement 11**). In principle, variation between

10   cells could be due to differential GpC MTase enzyme activity. However, the genome-wide levels of

11   GpC methylation reached comparable levels in all cells and the variability between cells was not

12   equally distributed across all DHS (**Fig. 2 d, Figure 2 –figure supplement 1**)

13   Measuring similarity of chromatin accessibility between cells was sufficient to group GM12878 and

14   K562 cells based on their cell type of origin (**Fig. 3 a**). In this particular case, the separation is

15   confounded with experimental batches. However, higher average GpC methylation in DHSs for the

16   respective cell type compared to the DHSs of the other cell type indicated that scNOMe-seq can

17   differentiate the two cell types (**Supplemental Fig. 14**). Similarly, endogenous CpG methylation at

18   different genomic features (DHS, 10 kb windows, gene bodies) was sufficient to distinguish

19   between the two cell types. This approach should be extendable to scNOMe-seq data from samples

20   containing mixtures of cell types.

21   scNOMe-seq measures chromatin accessibility at GpC di-nucleotides along the entire length of a

22   sequencing read. Since most features that bind DNA are smaller than the length of 100 bp (200 bp

23   within 200-50bp regions in the case of paired end reads), the regions covered by sequence-specific

24   transcription factors and nucleosomes can be captured within a single fragment. This allows one to

25   directly detect binding of TFs provided that their sequencing motif contains at least one GpC di-

nucleotide. I demonstrated the feasibility of this approach using CTCF binding sites. Of note, most motifs within regions of CTCF ChIP-seq peaks were protected from GpC methylation ('footprint') (**Fig. 5**). In agreement with an inferred binding event as the cause for this protection, scores for CTCF motifs that were associated with a footprint were significantly higher than for motifs without a footprint. Depending on the motif specificity of a given TF and provided that their motifs contain a GpC dinucleotide, similar measurements should be feasible for many TFs and could be used to infer the activity of a range of transcription factors in single cells or to measure combinatorial binding of two or more TFs.

Estimation of the average nucleosome phasing distances allows one to study chromatin compaction and complements the measurements of chromatin accessibility at regulatory regions and DNA methylation. The estimates from individual cells fit very well with measurements made from MNase-seq data in bulk samples(Valouev et al. 2011). It remains to be established whether the variation in phasing distances between individual cells is of biological or technical nature (**Fig. 6 e**).

These proof-of-principle experiments have been performed using commercial kits for bisulfite conversion and library amplification, additional optimization or alternative amplification approaches (Smallwood et al. 2014)are likely to increase the yield substantially. Compared to other single cell methods, for example ATAC-seq, scNOMe-seq does not enrich for accessible chromatin regions and thus requires significantly more sequencing coverage. Ultimately, it should be possible to integrate the GpC MTase treatment into microfluidic workflows and combine this method with scRNA-seq, similar to recently published methods that combine scRNA-seq and methylome-sequencing (Angermueller et al. 2016). This study was primarily designed to test the feasibility of NOMe-seq in single cells and only a small number of nuclei where sequenced for each cell line. As a consequence, this set up could not be used to study cell-to-cell variation in detail. scNOMe-seq will be particularly useful for studies that aim to simultaneously measure chromatin accessibility and DNA methylation. This approach will be especially powerful for the characterization of

12

1    chromatin organization in single cells from heterogeneous mixtures or complex tissues, for example

2    to samples of brain tissues or primary cancer cells.

3

## Methods

### Cell culture, nuclei isolation, and GpC methylase treatment

6    GM12878 (RRID:CVCL_7526) and K562 (RRID:CVCL_0004) cells were obtained directly from

7    Coriell and ATCC, respectively. No further confirmation of the authenticity of these cell lines or

8    mycoplasma testing has been performed. GM12878 were grown in RPMI medium 1640 (Gibco),

9    supplemented with 2mM L-Glutamine (Gibco), and Penicilin and Streptavidin (Pen Strep, Gibco),

10   and 15% fetal bovine serum (FBS, Gibco). K562 were grown in RPMI medium 1640 of the same

11   composition but with 10% FBS. Cells were grown at 37 C and in 5% $CO_2$. NOMe-Seq procedure

12   was performed based on protocols for CpG methyltransferase M.SSsI described in (Miranda et al.

13   2010) and the GpC methyltransferase from *M.CviPI* (Kelly et al. 2012), with some modification.

14   Between $2x10^6$ and $5x10^6$ cells were harvested by centrifuging the cell suspension for 5 min at

15   500x g. Cells were washed once with 1x PBS, re-suspended in 1 ml lysis buffer (10mM Tris-HCl

16   pH 7.4, 10mM NaCl, 3mM $MgCl_2$) and incubated for 10 min on ice. IGEPAL CA-630 (Sigma) was

17   added to a final concentration of 0.025% and the cell suspension was transferred to a 2 ml Dounce

18   homogenizer. Nuclei were released by 15 strokes with the pestle. Success of lysis was confirmed by

19   inspection under a light microscope. Nuclei were collected by centrifuging the cell suspension for 5

20   min at 800x g at 4 C and washed twice with cold lysis buffer without detergent. One million nuclei

21   were resupended in reaction buffer to yield a suspension with a final concentration of 1x GpC

22   MTase buffer (NEB), 0.32 mM S-Adenosylmethionine (SAM) (NEB), and 50 ul of GpC

23   methyltransferase (4U/ul)) from *M.CviPI* (NEB). The final reaction volume was 150 ul. The

24   suspension was carefully mixed before incubating for 8 min at 37 C after which another 25 ul of

25   enzyme and 0.7 ul of 32 mM SAM were added for an additional 8 min incubation at 37C. To avoid

13

1    disruption of nuclei incubation was stopped by adding 750 ul of 1x PBS and collecting the nuclei at

2    800 xg. Supernatant was removed and nuclei were re-suspended in 500ul 1x PBS containing

3    Hoechst 33342 DNA dye (NucBlue Live reagent, Hoechst). Nuclei were kept on ice until sorting.

4    For preparation of bulk libraries in GM21878 cell, nuclei preparation and GpC MTase treatment

5    was performed as described above. Nuclei were lysed immediately after incubation and DNA was

6    isolated using Phenol/Chloroform purification.

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

8    **conversion**

9    Nuclei were sorted at the Flow Cytometry core at the University of Chicago on a BD FACSAria or

10    BD FACSAria Fusio equipped with a 96-well-plate holder. To obtain individual and intact nuclei

11    gates were set on forward and side scatter to exclude aggregates and debris. DAPI/PacBlue channel

12    or Violet 450/500 channel were usedto excited the Hoechst 33342 DNA dye and to gate on cells

13    with DNA content corresponding to cells in G1 phase of the cell cycle in order to maintain similar

14    DNA content per cell and to remove potential heterogeneity attributable to cell cycle. Cells were

15    sorted into individual wells pre-filled with 19 ul of 1x M-Digestion buffer (EZ DNA Methylation

16    Direct Kit, Zymo Research) containing 1 mg/ml Proteinase K. Following collection, the plates were

17    briefly spun to collect droplets that might formed during handling. Nuclei were lysed by incubating

18    the samples at 50 C for 20 min in a PCR cycler. DNA was subjected to bisulfite conversion by

19    adding 130 ul of freshly prepared CT Conversion reagent (EZ DNA Methylation Direct Kit, Zymo)

20    to the lysed nuclei. Conversion was performed by denaturing the DNA at 98 C for 8 min followed

21    by 3.5 hrs incubation at 65 C. DNA isolation was performed using the EZ DNA Methylation Direct

22    Kit (Zymo Research) following the manufacturer's instruction with the modification that the DNA

23    was eluted in only 8 ul of elution buffer.

24    **Library preparation and sequencing**

14

1    Libraries were prepared using the Pico Methyl-seq Library prep Kit (Zymo Research) following the

2    manufacturer's instruction for low input samples. Specifically, the random primers were diluted 1:2

3    before the initial pre-amplification step and the first amplification was extended to a total of 10

4    amplification cycles. Libraries were amplified with barcoded primers allowing for multiplexing.

5    The sequences can be found in **Supplementary file 1**, primers were ordered from IDT. The

6    purification of amplified libraries was performed using Agencourt AMPureXP beads (Beckmann

7    Coulter), using a 1:1 ratio of beads and libraries. Concentration and size distribution of the final

8    libraries was assessed on an Bioanalyzer (Agilent). Libraries with average fragment size above 150

9    bp were pooled and sequenced. Libraries were sequenced on Illumina HiSeq 2500 in rapid mode

10    (K562 cells) and HiSeq4000 (GM12878 cells).

11    **Read processing and alignment**

12    Sequences were obtained using 100 bp paired-end mode. For processing and alignment each read

13    from a read pair was treated independently as this slightly improved the mapping efficiency. Before

14    alignment, read sequences in fastq format were assessed for quality using fastqc

15    (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). Reads were trimmed to remove low

16    quality bases and 6 bp were clipped from the 5 prime end of each read to avoid mismatches

17    introduced by amplification. In the case of GM12878 cells 6 bp were clipped from either end of the

18    read. Only reads that remained longer than 20 bp were kept for further analyses. These processing

19    steps were performed using trim_galore version 0.4.0

20    (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) with the following settings:

21    *trim_galore --quality 30 --phred33 --illumina --stringency 1 -e 0.1 --clip_R1 6 --gzip --length 20 --*

22    *output_dir outdir Sample.fastq.gz*. The trimmed fastq files were aligned using the bisulfite aligner

23    bismarck version 0.15.0 (Krueger et al. 2012) which calls bowtie2 (Langmead and Salzberg 2012)

24    internally. Reads were aligned to the human genome (genome assembly hg38). Reads were aligned

25    in single read mode using default settings. The amplification protocol used to generate the

26    scNOMe-seq libraries yielded non-directional libraries and alignment was performed with the

27    option —non_directional (*bismark --fastq --prefix SamplePrefix --output_dir output_dir --*

1    *non_directional --phred33-quals --score_min L,0,-0.2 --bowtie2 genome_file trimmed.fastq.gz*).

2    Some libraries contained small amounts of DNA from *C. elegans* as spike-ins, however these were

3    not used during the analysis. Duplicates were removed using samtools version 0.1.19 (Li et al.

4    2009) on sorted output files from bismark (*samtools rmdup SamplePrefix.sorted.bam*

5    *SampleAligned_rmdup.bam*).

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

7    Coverage and methylation status of all cytosines was extracted using

8    bismark_methylation_extractor (Krueger et al. 2012) (*bismark_methylation_extractor -s --ignore 6*

9    *--output outdir --cytosine_report --CX --genome_folder path_to_genome_data*

10   *SampleAligned_rmdup.bam*). The resulting coverage files were used to extract the methylation

11   status of cytosines specifically in GpC and CpG di-nucleotides using the coverage2cytosine script

12   which is part of Bismark (Krueger et al. 2012). The resulting coverage files contained cytosines in

13   GCG context which are ambiguous given that they represent a cytosine both in GpC and CpG di-

14   nucleotides. Coordinates of these ambiguous positions were identified using oligoMatch (Kent et al.

15   2002) and these positions were removed from the coverage files. The number of unconverted

16   cytosines (estimated based on apparent methylation rates in non-GpC and non-CpG context) was

17   low in all libraries (<1%). However, it was noted that unconverted cytosines were not randomly

18   distributed but associated with entirely unconverted reads. Regions covered by a read with more

19   than 3 unconverted cytosines in non-CpG and non-GpC context were removed from further analysis

20   as well. The genotype was not taken into account as its effect on calling the methylation status

21   incorrectly was deemed negligible for the analyses performed here.

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

23   ScNOMe-seq data were compared to a number of genomic features in GM12878 and K562 cells

24   collected by Encode (ENCODE Project ConsortiumThe ENCODE Project Consortium 2012) which

25   were downloaded through the UCSC data repository (Karolchik et al. 2014). These datasets are

16

1    listed in **Supplementary file 1**. While the scNOMe-seq data were aligned against human genome

2    assembly hg38, some of the datasets were only available on genome assembly hg19 and the

3    coordinates of these datasets were lifted from hg19 to hg38 using liftOver (Kent et al. 2002)

4    (default re-mapping ratio 1). Nucleosome positions based on MNase-seq data in GM12878 were

5    determined with DANPOS version 2.2.2 (Chen et al. 2013) using default settings. Resulting

6    intervals were lifted to hg38. After removing summit locations with occupancy values above 300,

7    the top 5% (713361) of nucleosome positions based on their summit occupancy value were used.

8    GpC and CpG methylation density across intervals encompassing DNase hypersensitivity sites

9    (DHSs), transcription factor binding sites (TFBS), and well positioned nucleosomes was calculated

10   across the 2 kb regions centered on the middle of these regions using the scoreMatrixBin function in

11   the genomation package (Akalin et al. 2015) in R (R Core Team 2015). Data were aggregated in 5

12   bp bins for each region and across all regions covered in a single cell. The average methylation

13   level in pre-defined intervals (DHSs, TFBS) was determined by computing the average GpC or

14   CpG methylation for each interval together with the number of GpC/CpGs covered in this interval

15   using the map function in bedtools (Quinlan and Hall 2010). If no other cut-offs were given, DHSs

16   were considered 'covered' and used in analyses when at least 4 GpCs occurring within the

17   predefined interval were covered by sequencing data in an individual cell. Because the frequency of

18   CpG di-nucleotides is significantly lower, only 2 CpGs were required in order for a DHSs to be

19   considered covered for analyses that focused on endogenous DNA methylation. To count the

20   number of cytosines within the primary sequence of a given DHSs only cytosines on the forward

21   strand were counted. While each GpC dinucleotide can be measured on both strands and would

22   therefore yield a count of two cytosines the data are sparse and each location will get at most a

23   single read. This approach should therefore give a more conservative estimate of the possible GpC

24   coverage. For analyses that used the scores of the peak regions, the peak scores reported the datasets

25   from bulk samples were used (ENCODE Project ConsortiumThe ENCODE Project Consortium

26   2012).

1 For analyses that were centered on transcription factor binding motifs the PWMs were obtained

2 from the JASPAR database (2014) (Tan) for the TFs CTCF (MA0139), EBF1 (MA0154), and

3 PU.1(MA0080). Genome-wide scanning for locations of sequence matches to the PWMs was

4 performed using matchPWM in the Biotstring package (Pages et al. 2016) in R with a threshold of

5 75% based on the human genome assembly hg38.

6 All plots were prepared using ggplot2 (Wickham 2009), with the exception of heatmaps displaying

7 the average methylation density around genomic features in individual cells which were prepared

8 using heatmap.2 in gplots (Warnes et al. 2016).

9 **Comparison of chromatin accessibility between cells**

10 Similarity in accessible chromatin between cells was calculated based on Jaccard similarity. Jaccard

11 similarity index (eq. 1) was calculated between pairs of samples by first obtaining the intersection

12 of DHSs covered in both samples of a pair with more than 4 GpCs. Each feature was annotated as

13 open or closed, depending on the methylation status (>= 40% methylation) and only pairs in which

14 at least one of the members was open were considered for this comparison.

$$jac(A,B) = \frac{(A \cap B)}{(A \cup B)} \tag{1}$$

15

16 The similarity between samples from GM12878 and K562 cells was calculated based on the union

17 set of DHSs from both cell lines. The similarity indexes of all pairwise comparisons were used to

18 compute the distances between each cell. The resulting clustered data were displayed as a heat map.

19 **CTCF footprints in single cells**

20 CTCF footprints were measured by comparing the GpC methylation level in each motif to the

21 methylation level in the 50bp flanking regions immediately upstream and downstream of the motif.

22 Overlapping motifs were merged into a single interval before determining the coordinates for

23 flanking regions. To ensure sufficient GpC coverage for each interval the resulting three adjacent

24 intervals for each locus were required to contain at least one covered GpC each, and 4 covered

25 GpCs in total. This analysis only included regions that were accessible based on the methylation

18

1    status of the flanking regions (at least 50%). A CTCF footprint 'score' was determined by simply

2    subtracting the average GpC methylation of the flanking regions from the GpC methylation of the

3    motif.

4    scNOMe-seq data were displayed in the UCSC genome browser (Kent et al. 2002) by converting

5    the GpC methylation coverage file into a bed file and using the methylation value as score. To

6    facilitated the visualization of the data in the context of previous Encode data the methylation files

7    were lifted to hg19. The tracks shown together with scNOMe-seq data are Open Chromatin by

8    DNaseI HS from ENCODE/OpenChrom (Duke University) for DNaseI hypersensitivity,

9    Nucleosome Signal from ENCODE/Stanford/BYU, and CTCF ChIP-seq signal from Broad Histone

10    Modification by ChIP-seq from ENCODE/Broad Institute. All data are from GM12878 cells.

11    **Estimation of nucleosome phasing**

12    Nucleosome phasing estimates were obtained by first calculating the correlation coefficients for the

13    methylation status of pairs of GpCs ad different offset distances. These values were computed using

14    a custom python script. Essentially, pairs of sequenced cytosines in GpC di-nucleotides were

15    collected for each offset distance from 3bp to 400bp cytosine. At each offset distance the correlation

16    of the methylation status was calculated across all pairs. Correlation coefficients were plotted

17    against the offset distances revealing periodic changes in the correlation coefficient. The

18    smoothened data were used to estimate the phasing distances by obtaining the offset distance

19    corresponding to the local maximum found between 100 bp and 300 bp. To determine phase lengths

20    of nucleosomes in different chromatin contexts the GpC coverage files were filtered for positions

21    falling into categories defined by chromHMM (ENCODE Project ConsortiumThe ENCODE

22    Project Consortium 2012; Ernst et al. 2011) before obtaining the correlation coefficients.

23    **Data access**

1   Raw data and methylation coverage files are available at GEO (https://www.ncbi.nlm.nih.gov/geo/)

2   under the accession number GSE83882. Reviewers might use this link:

3   http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=glotcwqqjbqlvef&acc=GSE83882

4

## Competing financial interest

6   The author declares no competing financial interests

## Acknowledgements

14

## References

16
17  Akalin A, Franke V, Vlahoviček K, Mason CE, Schübeler D. 2015. Genomation: a toolkit to
18      summarize, annotate and visualize genomic intervals. *Bioinformatics* **31**: 1127–1129.

19  Angermueller C, Clark SJ, Lee HJ, Macaulay IC, Teng MJ, Hu TX, Krueger F, Smallwood SEBA,
20      Ponting CP, Voet T, et al. 2016. Parallel single-cell sequencing links transcriptional and
21      epigenetic heterogeneity. *Nat Meth* 1–6.

22  Buenrostro JD, Wu B, Chang HY, Greenleaf WJ. 2015a. ATAC-seq: A Method for Assaying
23      Chromatin Accessibility Genome-Wide. *Current protocols in molecular biology / edited by
24      Frederick M Ausubel  [et al]* **109**: 21.29.1–9.

25  Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, Chang HY, Greenleaf
26      WJ. 2015b. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*
27      1–15.

1  Chen K, Xi Y, Pan X, Li Z, Kaestner K, Tyler J, Dent S, He X, Li W. 2013. DANPOS: dynamic
2      analysis of nucleosome position and occupancy by sequencing. *Genome Research* **23**: 341–351.

3  Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, Steemers FJ,
4      Trapnell C, Shendure J. 2015. Multiplex single cell profiling of chromatin accessibility by
5      combinatorial cellular indexing. *Science* **348**: 910–914.

6  ENCODE Project Consortium, The ENCODE Project Consortium. 2012. An integrated
7      encyclopedia of DNA elements in the human genome. *Nature* **489**: 57–74.

8  Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, Epstein CB, Zhang X, Wang L, Issner
9      R, Coyne M, et al. 2011. Mapping and analysis of chromatin state dynamics in nine human cell
10     types. *Nature* **473**: 43–49.

11 Farlik M, Sheffield NC, Nuzzo A, Datlinger P, Schönegger A, Klughammer J, Bock C. 2015.
12     Single-Cell DNA Methylome Sequencing and Bioinformatic Inference of Epigenomic Cell-
13     State Dynamics. *CellReports* **10**: 1386–1397.

14 Jaitin DA, Kenigsberg E, Keren-Shaul H, Elefant N, Paul F, Zaretsky I, Mildner A, Cohen N, Jung
15     S, Tanay A, et al. 2014. Massively parallel single-cell RNA-seq for marker-free decomposition
16     of tissues into cell types. *Science* **343**: 776–779.

17 Jin W, Tang Q, Wan M, Cui K, Zhang Y, Ren G, Ni B, Sklar J, Przytycka TM, Childs R, et al.
18     2015. Genome-wide detection of DNase I hypersensitive sites in single cells and FFPE tissue
19     samples. *Nature* 1–17.

20 Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, Dreszer TR, Fujita PA,
21     Guruvadoo L, Haeussler M, et al. 2014. The UCSC Genome Browser database: 2014 update.
22     *Nucleic Acids Research* **42**: D764–70.

23 Kelly TK, Liu Y, Lay FD, Liang G, Berman BP, Jones PA. 2012. Genome-wide mapping of
24     nucleosome positioning and DNA methylation within individual DNA molecules. *Genome
25     Research* **22**: 2497–2506.

26 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. 2002. The
27     human genome browser at UCSC. *Genome Research* **12**: 996–1006.

28 Kilgore JA, Hoose SA, Gustafson TL, Porter W, Kladde MP. 2007. Single-molecule and population
29     probing of chromatin structure using DNA methyltransferases. *Methods* **41**: 320–332.

30 Krueger F, Kreck B, Franke A, Andrews SR. 2012. DNA methylome analysis using short bisulfite
31     sequencing data. *Nat Meth* **9**: 145–151.

32 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Meth* **9**: 357–359.

33 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
34     1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format
35     and SAMtools. *Bioinformatics* **25**: 2078–2079.

36 Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, Tirosh I, Bialas AR, Kamitaki
37     N, Martersteck EM, et al. 2015. Highly Parallel Genome-wide Expression Profiling of
38     Individual Cells Using Nanoliter Droplets. *Cell* **161**: 1202–1214.

39 Maurano MT, Stamatoyannopoulos JA. 2015. Taking Stock of Regulatory Variation. *Cell Systems*

**1**: 18–21.

Miranda TB, Kelly TK, Bouazoune K, Jones PA. 2010. Methylation-sensitive single-molecule analysis of chromatin structure. *Current protocols in molecular biology / edited by Frederick M Ausubel [et al]* **Chapter 21**: Unit 21.17.1–16.

Neph S, Vierstra J, Stergachis AB, Reynolds AP, Haugen E, Vernot B, Thurman RE, John S, Sandstrom R, Johnson AK, et al. 2012. An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489**: 83–90.

Pages H, Aboyoun P, Gentleman RC, DebRoy S. 2016. *Biostrings: String objects representing biological sequences, and matching algorithms*.

Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**: 841–842.

R Core Team. 2015. *R: A language and environment for statistical computing*. https://www.R-project.org/.

Rotem A, Ram O, Shoresh N, Sperling RA, Goren A, Weitz DA, Bernstein BE. 2015. Single-cell ChIP-seq reveals cell subpopulations defined by chromatin state. *Nature Biotechnology* **33**: 1–11.

Schones DE, Cui K, Cuddapah S, Roh T-Y, Barski A, Wang Z, Wei G, Zhao K. 2008. Dynamic Regulation of Nucleosome Positioning in the Human Genome. *Cell* **132**: 887–898.

Schübeler D. 2015. Function and information content of DNA methylation. *Nature* **517**: 321–326.

Shalek AK, Satija R, Shuga J, Trombetta JJ, Gennert D, Lu D, Chen P, Gertner RS, Gaublomme JT, Yosef N, et al. 2014. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* **510**: 263–269.

Small EC, Xi L, Wang J-P, Widom J, Licht JD. 2014. Single-cell nucleosome mapping reveals the molecular basis of gene expression heterogeneity. *Proc Natl Acad Sci USA* **111**: E2462–71.

Smallwood SA, Lee HJ, Angermueller C, Krueger F, Saadeh H, Peat J, Andrews SR, Stegle O, Reik W, Kelsey G. 2014. Single-cell genome-wide bisulfite sequencing for assessing epigenetic heterogeneity. *Nat Meth* **11**: 817–820.

Stadler MB, Murr R, Burger L, Ivanek R, Lienert F, Schöler A, van Nimwegen E, Wirbelauer C, Oakeley EJ, Gaidatzis D, et al. 2011. DNA-binding factors shape the mouse methylome at distal regulatory regions. *Nature* 1–7.

Taberlay PC, Statham AL, Kelly TK, Clark SJ, Jones PA. 2014. Reconfiguration of nucleosome-depleted regions at distal regulatory elements accompanies DNA methylation of enhancers and insulators in cancer. *Genome Research* **24**: 1421–1432.

Tan G. JASPAR2014: Data package for JASPAR. http://jaspar.genereg.net/.

Treutlein B, Brownfield DG, Wu AR, Neff NF, Mantalas GL, Espinoza FH, Desai TJ, Krasnow MA, Quake SR. 2014. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509**: 371–375.

Valouev A, Johnson SM, Boyd SD, Smith CL, Fire AZ, Sidow A. 2011. Determinants of

22

1    nucleosome organization in primary human cells. *Nature* **474**: 516–520.

2    Warnes GR, Bolker B, Bonebakker L, Gentleman R. 2016. *gplots: Various R programming tools*
3        *for plotting data*. R package version  https://CRAN.R-project.org/package=gplots.

4    Wickham H. 2009. *ggplot2*. Springer Science & Business Media, New York, NY.

5    Ziller MJ, Edri R, Yaffe Y, Donaghey J, Pop R, Mallard W, Issner R, Gifford CA, Goren A, Xing J,
6        et al. 2014. Dissecting neural differentiation regulatory networks through epigenetic
7        footprinting. *Nature* 1–16.

8

9
10
11
12    **Figure Legends**

13    **Figure 1 scNOMe-seq detected DNase hypersensitive sites in single cells**. a) Schematic of GpC

14    methyltransferase-based mapping of chromatin accessibility and simultaneous detection of

15    endogenous DNA methylation. b) Schematic of scNOMe-seq procedure introduced in this study.

16    **Figure 1 – figure supplement 1: FACS profile from Hoechst stained nuclei to assess DNA**

17    **content.** Nuclei were stained with Hoechst 33342 DNA dye and nuclei with DNA content

18    corresponding to the G1-phase of the cell cycle were sorted into individual wells in a 96 well plate.

19    Aggregates and debris were removed using gates on forward and side scatter.

20

21    **Figure 1 – figure supplement 2: Schematic of experimental set up**. A total of 19 individual cells

22    from GM12878 were profiled in this study, 12 of these cells were exposed to GpC MTase and 7

23    were subjected to the same process without exposure to MTase. For K562 11 cells were profiled all

24    of which were subjected to GpC MTase treatment.

25

26    **Figure 1 – figure supplement 3: Number of covered GpC and CpG dinucleotides is**

27    **proportional to the number of total bases covered**. Number of covered cytosines in GpC and

28    CpG dinucleotides plotted against the total number of nucleotides covered per sample. This

29    comparison suggests that there is no strong bias towards or against GpC and CpG dinucleotides.

30    This plot also shows that the coverage was about 2-fold higher for K562 cells compared to

31    GM12878.

32

33

34    **Figure 2 scNOMe-seq data reveal how accessibility in single cells underlies observed DNaseI**

35    **hypersensitivity in a population of cells**. a) Average GpC methylation level (blue) and CpG

23

1    methylation level (orange) at DHSs in GM12878 cells. Regions are centered on the middle of

2    DNase-seq peak locations. Shown is the average methylation across a 2 kb window of 12 GM12878

3    cells. b) Heatmap displaying the average GpC methylation level across the same regions as in a).

4    Each row corresponds to an individual GM12878 cell. Cells were grouped by similarity. c)

5    Proportion of DHSs covered by scNOMe-seq sequencing reads in each cell. The proportion

6    displayed corresponds to the fraction of DHSs covered by at least 1 or 4 GpCs in a given cell. Only

7    DHSs with at least 1 GpC (red) or 4 GpCs (cyan) within their primary sequence were taken in

8    consideration. Error bars represent standard deviation. d) Average GpC methylation at DHSs

9    grouped into quartiles based on associated DNase-seq peak scores from lowest to highest scores.

10   'Shuffled' represents methylation data in genomic regions obtained by random placements of DHS

11   peak intervals. Data shown are from GM12878 cells. e) Fraction of accessible sites in individual

12   GM12878 cells (red) and K562 cells (cyan). Shown are the means and standard deviation based on

13   all cells. f) Scatter plot showing relationship between GpC methylation levels and DHS peaks score

14   for each covered DHS. Plot shows data from all individual GM12878 cells. Red trend line is shown

15   to visualize the relationship between GpC methylation and endogenous CpG methylation. g) Scatter

16   plot showing relationship between CpG methylation levels and DHS peaks score for each covered

17   DHS. Plot shows data from all individual GM12878 cells. Red trend line is shown to visualize the

18   relationship between CpG methylation and peak scores. h) Plot illustrates the relationship between

19   endogenous CpG methylation and GpC methylation at DHS loci. Plot shows combined data from

20   all GM12878 cells. Correlation was calculated based on Pearson correlation (r = -0.13) i) Average

21   CpG methylation at DHS loci grouped based on GpC scores within single cells. Each dot represents

22   the average CpG methylation level for a single cell.

23

24   **Figure 2 – figure supplement 1: Average CpG and GpC methylation levels in single cells**.

25   Boxplots representing the methylation level at CpG and GpC dinucleotides for groups of cells

26   (GM12878 w/ and w/o MTase,K562 w/ MTase). GM12878 and K562 cells show different levels of

27   CpG methylation. The difference in CpG methylation between GM12878 w/o MTase and

28   GM12878 w/ MTase treatment was largely driven by two cells. These cells were kept as no other

29   criterion suggested their removal. GpC MTase treated cells shows a clear enrichment of GpC

30   methylation while GM12878 cells not exposed to MTase do not show levels above 1%. These

31   might reflect incomplete conversion, minimal cross-contamination during the parallel preparation,

32   or activity of endogenous methyltransferases.

33

34   **Figure 2 – figure supplement 2: Heatmaps of average GpC and CpG methylation across DHS**

35   **regions in GM12878 cells.** Each row represents data from an individual cell, both treated and

1 control samples are plotted together. Cells were grouped using hierarchical clustering based on GpC

2 methylation (left) and CpG methylation (right) within 2 kb regions around DHSs. As expected GpC

3 methylation clearly separates MTase treated and untreated samples. Endogenous CpG methylation

4 does not differ systematically between MTase treated and untreated samples.

5

6 **Figure 2 – figure supplement 3: Average GpC and CpG methylation across DHS regions in**

7 **K562 cells**.

8 Average GpC methylation level (blue) and CpG methylation level (orange) at DNase

9 Hypersensitive sites (DHSs) in K562 cells. Regions are centered on the middle of DNase-seq peak

10 locations. Shown is the average methylation across a 2 kb window of the pool of 11 K562 cells.

11 **Figure 2 – figure supplement 4: Heatmaps of average GpC and CpG methylation across DHS**

12 **regions in K562 cells.** Each row represents data from an individual cell. Cells were grouped using

13 hierarchical clustering based on GpC methylation (left) and CpG methylation (right) within 2kb

14 regions around DHSs.

15

16 **Figure 2 – figure supplement 5: Distribution of counts of GpCs within DHSs in GM12878 and**

17 **K562 cells.** Histogram shows the number of GpCs per DHS in GM12878 cells (left) and K562 cells

18 (right). While each GpC dinucleotide can be measured on both strands and would therefore yield a

19 count of two cytosines this histogram only displays counts per GpC (using the cytosines on the

20 forward strand). This is to account for the fact that in single cells DHSs will be covered at most by

21 one or two reads that originate from the same fragment.

22

23 **Figure 2 – figure supplement 6: Proportion of DHSs at different cutoffs for GpCs and CpGs.**

24 Bar graphs display proportion of DHSs that contain at least the number of GpCs (left) and CpGs

25 (right) indicated. The proportions are given in relation to the total number of DHSs in each cell line.

26 Numbers in the bars refer to the number of DHSs at that GpC and CpG threshold. As described in

27 **Figure 2 – figure supplement 5** and methods, only cytosines in GpC and CpG di–nucleotides on

28 the forward strand were counted for this analysis.

29

30 **Figure 2 – figure supplement 7: Relationship between DNase-seq peak score and GpC and**

31 **CpG methylation in GM12878 and K562 cells** a) Average GpC methylation and b), c)

32 endogenous CpG methylation at DHSs grouped into quartiles based on associated DNase-seq peak

33 scores from lowest to highest scores. 'Shuffled' represents methylation data in genomic regions

25

1   obtained by random placements of DHS peak intervals. a) and c) show data for 11 K562 cells and

2   b) shows data for 12 GM12878 cells. Each point represents average score from a single cell.

3   **Figure 2 – figure supplement 8: Correlation between GpC methylation and DHS peak score**

4   Shown are correlation coefficients for comparisons between single cell and bulk NOMe-seq data

5   with DNase-seq peak score for each covered location for a) GM12878 and b) K562. Each dot

6   represents value for a single Pearson correlation. The correlation between GpC methylation and

7   DHS peaks scores was significantly lower in single cells compared to bulk NOMe-seq data. No

8   correlation was observed between GpC methylation and DHS peak score using randomized DHS

9   locations.

10  **Figure 2 – figure supplement 9: Cumulative distribution of average GpC methylation in DHSs**

11  **in GM12878 and K562 cells.** Plot of cumulative distribution of GpC methylation for individual

12  GM12878 and K562 cells at DHSs with at least 4 covered GpC. GM12878 and K562 cells exposed

13  to GpC MTase show similar distributions. About 50% of all cells show no or low methylation (<=

14  25%). GM12878 cells not exposed to GpC MTase do not show any significant number of DHSs

15  with GpC methylation.

16

17  **Figure 2 – figure supplement 10: Proportion of accessible DHSs remains stable across range**

18  **of thresholds for methylation levels and covered GpCs per site.**

19  Different thresholds for GpC methylation and number of covered GpC required per individual DHS

20  were used to test how much the number of resulting 'accessible' DHSs depended on these

21  parameters. Threshold for GpC methylations was varied between 25% and 50% while the number

22  of required GpCs for a DHSs to be considered in this analysis was varied between 1 and 8. The

23  proportion of accessible sites is plotted for each set of parameters in GM12878 cells (left) and K562

24  cells (right). Proportion of accessible sites remained relatively stable across the range of parameters.

25  Note that the categories with low GpCs count thresholds contain all sites above this threshold.

26

27  **Figure 2 – figure supplement 11: Cell-to-cell variability in DHSs accessibility reflects DNaseI**

28  **hypersensitivity of the region.** Pair-wise jaccard distances between GM12878 (a) and K562 (b)

29  cells, respectively, were calculated based on DHSs accessibility in individual cells. DHSs were

30  grouped by DNase-seq peak scores and DHSs were considered accessible if the average

31  methylation for that locus was above 40%. Only DHSs with at least 4 covered GpCs were included

32  in this analysis.

33

26

1    **Figure 2 – figure supplement 12: Comparison of correlations between single cell NOMe-seq**

2    **and bulk NOMe-seq data sets.** Shown are the correlation coefficients for comparisons between

3    GpC methylation in single cell and bulk NOMe-seq data within DNase-seq peak location for a)

4    GM12878 and b) K562. Each dot represents value for a single Pearson correlation. Comparison was

5    performed on original DHS loci (left) and on randomized DHS loci (right).

6

7    **Figure 2 – figure supplement 13: GpC methylation correlates with DHS peaks scores in**

8    **individual cells.** High DHS peak scores are associated with higher GpC methylation in single cells.

9    Scatter plot showing relationship between GpC methylation levels and DHS peaks scores for each

10   covered DHS. Each plot shows data from an individual GM12878 cell. Red trend line to aid

11   visualization of the relationship between GpC methylation and peak scores.

12

13   **Figure 2 – figure supplement 14: Endogenous CpG methylation is inversely correlated with**

14   **DHS peak scores in individual cells.** High DHS peak scores are associated with lower endogenous

15   CpG methylation in single cells. Scatter plot showing relationship between CpG methylation levels

16   and DHS peaks score for each covered DHS. Each plot shows data from an individual GM12878

17   cell. Red trend line to aid visualization of the relationship between GpC methylation and peak

18   scores.

19

20   **Figure 2 – figure supplement 15: Comparison of CpG and GpC methylation status at**

21   **individual DHS in single GM12878 cells.** Smoothened scatterplot illustrates the relationship

22   between endogenous CpG methylation and GpC methylation at DHS loci. Each plot shows data

23   from a single GM12878 cell.

24

25

26   **Figure 3 Single cell NOMe-seq reveals chromatin features closely linked to gene expression.**

27   a) Average GpC methylation level at TSS in GM12878 cells. Regions are centered on the TSS

28   locations. Shown is the average methylation across a 2 kb window of 12 GM12878 cells. b) Same

29   as in a) but displaying the endogenous CpG methylation level. C) Average endogenous CpG

30   methylation at gene loci in individual GM12878 cells. Shown is the average methylation across

31   gene bodies (represented as meta genes) and 50 kb regions upstream and downstream of each gene.

32   Each line represents the aggregated CpG methylation data for a single GM12878 cell (TES:

33   transcription end site). d) Boxplot displays average CpG methylation in gene bodies. Genes were

34   grouped into quartiles based on their transcript levels in bulk. Dots represent the average CpG

27

1   methylation value for individual cells. e) Boxplot displays average GpC methylation in promoter

2   regions (-500 bp to +150 bp). Genes were grouped into quartiles based on their transcript levels in

3   bulk. f) Similar to e) but displayed are the levels of endogenous CpG methylation.

4   **Figure 3 – figure supplement 1: Endogenous methylation in gene bodies of single K562 cells** a)

5   Average endogenous CpG methylation at gene loci in individual K562 cells. Shown is the average

6   methylation across gene bodies (represented as meta genes) and 50 kb regions upstream and

7   downstream of each gene. Each line represents the aggregated CpG methylation data for a single

8   K562 cell (TES: transcription end site). b) Boxplot displays average CpG methylation in gene

9   bodies. Genes were grouped into quartiles based on their transcript levels in bulk. Dots represent the

10  average CpG methylation value for individual cells.

11  **Figure 3 – figure supplement 2: Chromatin accessibility in promoters correlates with**

12  **transcript levels of adjacent genes** a) Average GpC methylation level at TSS genes in GM12878

13  cells. Regions are centered on the TSS locations and genes were grouped into quartiles based on

14  their transcript levels in bulk GM12878 cells. b) The same plot as in a) based on scNOMe-seq data

15  from K562 cells and, correspondingly, transcript levels in K562 cells.

16

17  **Figure 4: single cell GpC and CpG methylation signal is sufficient to group GM12878 and**

18  **K562 cells according to their origin** a) Heatmap shows similarity scores (pair-wise Jaccard

19  distances) for accessibility between all GM12878 and K562 cells measured on the union set of

20  DHSs from GM12878 and K562 cells. Cells were grouped based on unsupervised hierarchical

21  clustering. b) Average GpC methylation at the DHSs from GM12878 cells and K562 cells,

22  respectively, was calculated for all individual GM12878 and K562 cells. The resulting two values

23  for GpC methylation are displayed for each cell. GM12878 and K562 are separable based on these

24  data. GM12878 and K562 cells showed different levels of genome-wide GpC methylation.

25  Consequently, the average methylation levels at K562 DHSs for both cell types are similar.

26  However, for cells from either cell type the methylation levels are higher in the DHSs of the cell

27  type of origin than in the DHSs of the other cell type. c) Heatmap shows correlation coefficients

28  between all GM12878 and K562 cells for pair-wise comparison of CpG methylation levels.

29  Genome was divided into 10 kb bins and only bins with sufficient coverage in both cells were used

30  for a given pair (>= 20 covered CpGs). Cells were grouped based on unsupervised hierarchical

31  clustering.

32  **Figure 4 – figure supplement 1 Single GM12878 and K562 cells can be grouped based on GpC**

33  **methylation and endogenous methylation** a) Heatmap shows correlation coefficients for GpC

28

1    methylation between all GM12878 and K562 cells measured on the union set of DHSs from

2    GM12878 and K562 cells. Cells were grouped based on unsupervised hierarchical clustering. Only

3    DHS with at least 4 covered GpCs in both cells were used for pair-wise comparison. b) Same as in

4    a) but based on CpG methylation level in the union set of DHSs, at least 2 covered CpGs in both

5    cells were required to include a DHS in the pair-wise comparison. c) Heatmap shows correlation

6    coefficients between all GM12878 and K562 cells for pair-wise comparison of CpG methylation

7    levels in gene bodies. Only loci with sufficient coverage in both cells were used for a given pair (>=

8    10 covered CpGs). Cells were grouped based on unsupervised hierarchical clustering. All

9    correlation coefficients were calculated using Pearson correlation.

10

11    **Figure 5**. **scNOMe-seq detected characteristic accessibility patterns at CTCF transcription**

12    **factor binding sites and measured CTCF footprints at individual loci** a) Average GpC

13    methylation level (blue) and CpG methylation level (orange) at CTCF binding sites in GM12878

14    cells. Regions are centered on motif locations. Shown is the average methylation across a 2 kb

15    window of the pool of 12 GM12878 cells. b) Heatmap displaying the average GpC methylation

16    across CTCF binding sites. Each row corresponds to an individual GM12878 cell and rows are

17    grouped by similarity. c) Schematic outline the measurement of CTCF footprints in accessible

18    regions. M denotes CTCF binding motifs within CTCF ChIP-seq regions and U and D indicate 50

19    bp upstream and downstream flanking regions. footprint score was determined by subtracting the

20    average GpC methylation in the flanking regions from the GpC methylation at the motif. d)

21    Heatmap displays GpC methylation in accessible regions found in a representative GM12878 cell

22    (GM_1). Each row represents a single CTCF motif instance within a CTCF ChIP-seq region.

23    Average methylation values for the motif and the 50 bp upstream and downstream regions are

24    shown separately. Regions are sorted based on the footprint score. Displayed are only regions that

25    had sufficient GpC coverage and that were considered accessible based on the methylation status of

26    the flanking regions. e) Heatmap reporting the CTCF motif scores for the motif regions in d).

27    Regions are sorted in the same order as in d). f) Average number of accessible regions at CTCF

28    motifs and the average number of those with a detectable footprint per individual GM12878 cell.

29    Error bars reflect standard deviation. g) Average CTCF motif scores in regions with and without

30    CTCF footprint for all 12 GM12878 cells. Each line connects the two data points from an individual

31    cell h) Combined display of scNOMe-seq data from this study and DNase hypersensitivity data,

32    nucleosome occupancy, and CTCF ChIP-seq data from ENCODE. Upper panel shows a ~10 kb

33    region containing a CTCF binding site. DNaseI hypersensitivity data and nucleosome density show

34    characteristic distribution around CTCF binding sites in GM12878 cells. Lower panel shows the

1  GpC methylation data of 5 individual cells that had sequencing coverage in this region, 4 of the

2  cells provide GpC data covering the CTCF motif located in the region. scNOMe-seq data tracks

3  show methylation status of individual GpCs. Each row corresponds to data from a single cell. These

4  data indicate that binding of CTCF is detected in all 4 cells. Data are displayed as tracks in the

5  UCSC genome browser (*http://genome.ucsc.edu)*.

6  **Figure 5 – figure supplement 1: Average GpC methylation and endogenous CpG methylation**

7  **at CTCF sites in pooled K562 cells** Average GpC methylation level (blue) and CpG methylation

8  level (orange) at CTCF binding sites in K562 cells. Regions are centered on motif locations. Shown

9  is the average methylation across a 2 kb window of the pool of 11 K562 cells.

10

11  **Figure 5 – figure supplement 2: Average GpC methylation level at CTCF binding sites in**

12  **individual K562 cells**. Heatmap shows the average GpC methylation across a 2 kb window

13  centered on the CTCF motif location. Each row corresponds to an individual K562 cell and rows

14  are grouped by hierarchical clustering.

15

16  **Figure 5 – figure supplement 3: Average GpC methylation and endogenous CpG methylation**

17  **at additional transcription factor binding sites in pools of GM12878 and K562 cells.** Average

18  GpC methylation level (blue) and CpG methylation level (orange) at a) PU.1 binding sites in

19  GM12878 cells and b) EBF1 binding sites. Regions are centered on motif locations. Shown is the

20  average methylation across a 2 kb window of the pool of 12 GM12878 cells. c) Same plot as in a)

21  but based on PU.1 binding sites in K562 cells. Regions are centered on motif locations. Shown is

22  the average methylation across a 2 kb window of the pool of 11 K562 cells.

23

24  **Figure 5 – figure supplement 4: Scores at CTCF motifs with footprints are significantly**

25  **higher than those without.** Boxplot representing the CTCF motif scores in regions with and

26  without CTCF footprint of an individual GM12878 cell (GM_1, the same cell as shown in Figure 3

27  d and e)).

28

29  **Figure 5 – figure supplement 5: Loci with CTCF footprint in single cells**. At each locus

30  scNOMe-seq data from this study and DNase hypersensitivity data from ENCODE are shown.

31  scNOMe-seq data tracks show methylation status of individual GpCs. Each row corresponds to data

32  from a single cell, Colors indicate the methylation status of each GpC (yellow: methylated; blue:

33  unmethylated). Data are displayed as tracks in the UCSC genome browser.

34

35

30

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

11  **Figure 6 – figure supplement 1: Number of nucleotide pairs used for correlation at each offset**
12  **distance**. Plotted is the number of nucleotide pairs that are found at each offset distance and used to
13  calculate the correlation coefficient at that distance. The number of comparison declines
14  precipitously. While each read has a maximum of 100 bp all samples were sequenced in paired-end
15  mode and even though the alignment was performed in single end mode additional comparisons can
16  therefore be made in most cases based on GpCs covered in the read from the opposite side of a
17  fragment. In many cases, the data become too sparse beyond 400 bp. Each line represents data from
18  an individual cell. Data from GM12878 and K562 cells are plotted on the left and right,
19  respectively.

20  **Figure 6 – figure supplement 2: Offset distances with a high proportion of GpC pairs that**
21  **share methylation status indicate average phasing distance.** Shown is the proportion of
22  nucleotide pairs at each offset distance in which both cytosines are methylated. These
23  measurements yield curves very similar to the distribution of Pearson correlation coefficients
24  (Figure 6 c and d).

25

26  **Supplementary file 1: Table 1:** scNOMe-seq libraries used in this paper and their technical details
27  and alignment summary statistics. **Table 2:** Primer sequences of primers used for amplification and
28  barcoding of sequencing library. **Table 3:** Additional datasets used in this study and their sources.

29
30

**Figure Legends**

**Figure 1: scNOMe-seq detected DNase hypersensitive sites in single cells**. a) Schematic of GpC methyltransferase-based mapping of chromatin accessibility and simultaneous detection of endogenous DNA methylation. b) Schematic of scNOMe-seq procedure introduced in this study.

**Figure 2: scNOMe-seq data reveal how accessibility in single cells underlies observed DNaseI hypersensitivity in a population of cells**. a) Average GpC methylation level (blue) and CpG methylation level (orange) at DHSs in GM12878 cells. Regions are centered on the middle of DNase-seq peak locations. Shown is the average methylation across a 2 kb window of 12 GM12878 cells. b) Heatmap displaying the average GpC methylation level across the same regions as in a). Each row corresponds to an individual GM12878 cell. Cells were grouped by similarity. c) Proportion of DHSs covered by scNOMe-seq sequencing reads in each cell. The proportion displayed corresponds to the fraction of DHSs covered by at least 1 or 4 GpCs in a given cell. Only DHSs with at least 1 GpC (red) or 4 GpCs (cyan) within their primary sequence were taken in consideration. Error bars represent standard deviation. d) Average GpC methylation at DHSs grouped into quartiles based on associated DNase-seq peak scores from lowest to highest scores. 'Shuffled' represents methylation data in genomic regions obtained by random placements of DHS peak intervals. Data shown are from GM12878 cells. e) Fraction of accessible sites in individual GM12878 cells (red) and K562 cells (cyan). Shown are the means and standard deviation based on all cells. f) Scatter plot showing relationship between GpC methylation levels and DHS peaks score for each covered DHS. Plot shows data from all individual GM12878 cells. Red trend line is shown to visualize the relationship between GpC methylation and endogenous CpG methylation. g) Scatter plot showing relationship between CpG methylation levels and DHS peaks score for each covered DHS. Plot shows data from all individual GM12878 cells. Red trend line is shown to visualize the relationship between CpG methylation and peak scores. h) Plot illustrates the relationship between endogenous CpG methylation and GpC methylation at DHS loci. Plot shows combined data from all GM12878 cells. Correlation was calculated based on Pearson correlation (r = -0.13) i) Average CpG methylation at DHS loci grouped based on GpC scores within single cells. Each dot represents the average CpG methylation level for a single cell.

**Figure 3: Single cell NOMe-seq reveals chromatin features closely linked to gene expression.**
a) Average GpC methylation level at TSS in GM12878 cells. Regions are centered on the TSS locations. Shown is the average methylation across a 2 kb window of 12 GM12878 cells. b) Same as in a) but displaying the endogenous CpG methylation level. C) Average endogenous CpG methylation at gene loci in individual GM12878 cells. Shown is the average methylation across gene bodies (represented as meta genes) and 50 kb regions upstream and downstream of each gene. Each line represents the aggregated CpG methylation data for a single GM12878 cell (TES: transcription end site). d) Boxplot displays average CpG methylation in gene bodies. Genes were grouped into quartiles based on their transcript levels in bulk. Dots represent the average CpG methylation value for individual cells. e) Boxplot displays average GpC methylation in promoter regions (-500 bp to +150 bp). Genes were grouped into quartiles based on their transcript levels in bulk. f) Similar to e) but displayed are the levels of endogenous CpG methylation.

**Figure 4: single cell GpC and CpG methylation signal is sufficient to group GM12878 and K562 cells according to their origin** a) Heatmap shows similarity scores (pair-wise Jaccard distances) for accessibility between all GM12878 and K562 cells measured on the union set of DHSs from GM12878 and K562 cells. Cells were grouped based on unsupervised hierarchical clustering. b) Average GpC methylation at the DHSs from GM12878 cells and K562 cells, respectively, was calculated for all individual GM12878 and K562 cells. The resulting two values for GpC methylation are displayed for each cell. GM12878 and K562 are separable based on these data. GM12878 and K562 cells showed different levels of genome-wide GpC methylation. Consequently, the average methylation levels at K562 DHSs for both cell types are similar. However, for cells from either cell type the methylation levels are higher in the DHSs of the cell type of origin than in the DHSs of the other cell type. c) Heatmap shows correlation coefficients between all GM12878 and K562 cells for pair-wise comparison of CpG methylation levels. Genome was divided into 10 kb bins and only bins with sufficient coverage in both cells were used for a given pair (>= 20 covered CpGs). Cells were grouped based on unsupervised hierarchical clustering.

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

**Figure 6: Nucleosome phasing in single cells**. a) Average GpC methylation level and b) CpG methylation level at well-positioned nucleosomes in GM12878 cells. Regions are centered on

midpoints of top 5% of positioned nucleosomes. Shown is the average methylation across a 2 kb window of the pool of 12 GM12878 cells. c), d) Correlation coefficients for the comparison in methylation status between GpCs separated by different offset distances for GM12878 (c) and K562 (d) cells. Each line represents a single cell. Data are smoothened for better visualization. e) Distribution of estimated phase lengths for GM12878 and K562 cells. f) Nucleosome phasing in GM12878 in genomic regions associated with different chromatin states defined by chromHMM (ENCODE). Boxplot represents the distribution of estimated phase lengths from all 12 GM12878 cells and overlaid points indicate values of each individual cells.
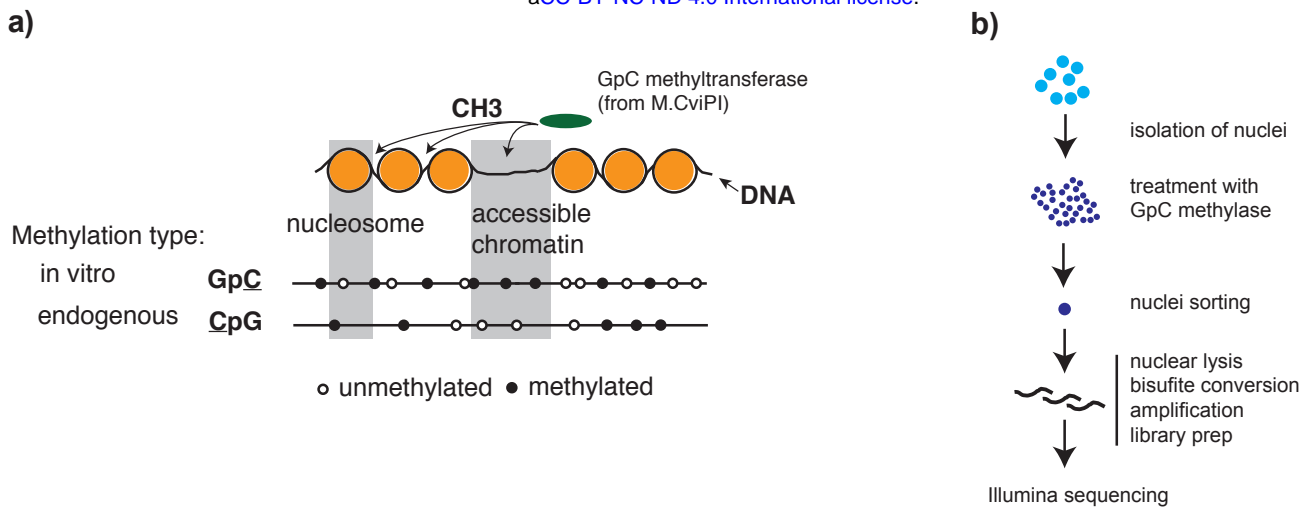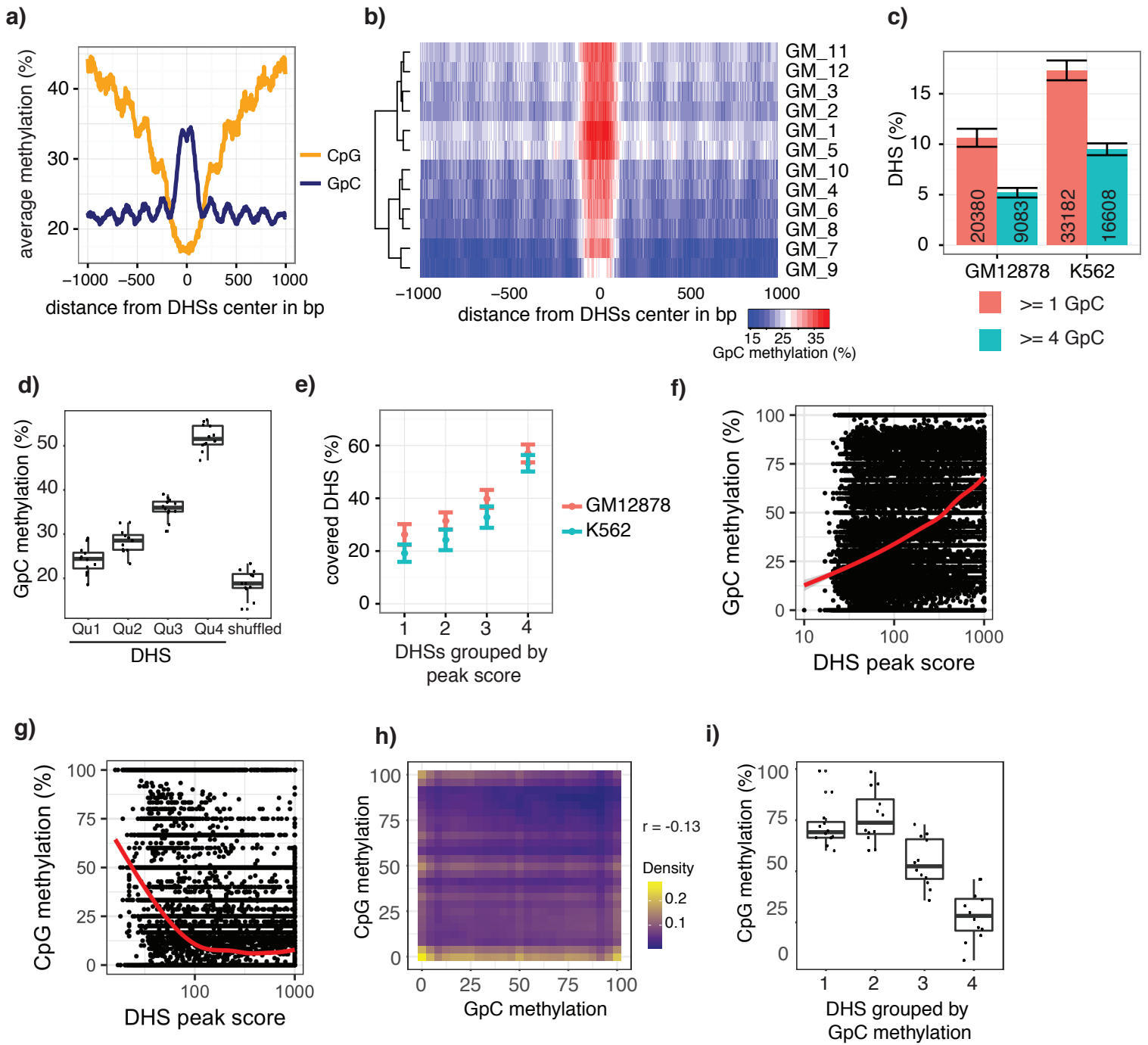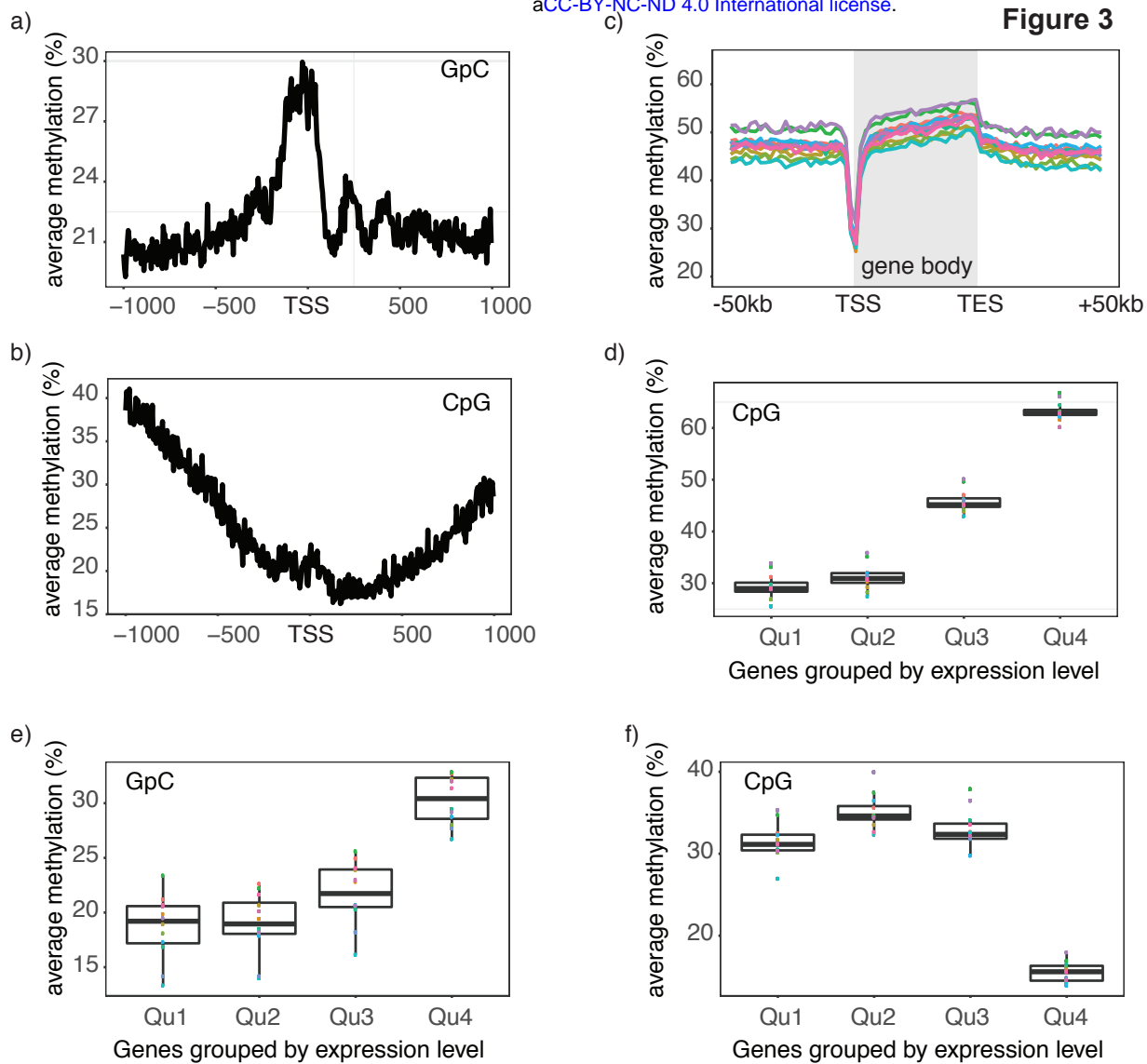
**Figure 1**

**a)**

GpC methyltransferase
(from M.CviPI)

**CH3**

**DNA**

nucleosome

accessible
chromatin

Methylation type:

in vitro **GpC**

endogenous **CpG**

○ unmethylated ● methylated

**b)**

isolation of nuclei

treatment with
GpC methylase

nuclei sorting

nuclear lysis
bisufite conversion
amplification
library prep

Illumina sequencing

**Figure 2**

**Figure 3**

**Figure 4**

**Figure 5**

**Figure 6**