1    High coverage of single cell genomes by T7-assisted enzymatic methyl-sequencing

2    Juan Wang[1,†], Yitong Fang[1,†], WenFang Chen[1,†], Chen Zhang[1],Zhichao Chen[1,2], Zhe Xie[2,3],Zhe
3    Weng[1], Weitian Chen[1,2],Fengying Ruan[1], , Yeming Xie[1], Yuxin Sun[1], Mei Guo[1], Yaning Li[1],
4    Chong Tang[1,4,5,*]
5

6    [1.] BGI. Genomics, BGI-Shenzhen, Shenzhen 518083, China
7    [2.] College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049, China
8    [3.] Department of Biology, Cell Biology and Physiology, University of Copenhagen 13, 2100
9    Copenhagen，Denmark

10    **†** These authors contributed equally to this work.

11
12    Keywords: single cell, methylation, CpG, coverage, EM-seq
13
14    *Running title: High single cell genome coverage by TEAM-seq*

15

16    *Correspondence:

17    Chong Tang

18    Director of Technology, BGI-Shenzhen, China

19    Phone: +8618025420976

20    Email: tangchong@bgi.com

21

1

**Abstract**

Conventional approaches to studying 5mC marks in single cells or samples with picogram input DNA amounts usually suffer from low genome coverage due to DNA degradation. Many methods have been developed to optimize the library construction efficiency for bisulfite-treated DNA. However, most of these approaches ignored the amplification bias of bisulfite-treated DNA, which leads to shallow genome coverage. In this study, we developed the T7-assisted enzymatic methyl-sequencing method (TEAM-seq), which adopts enzymatic conversion to minimize DNA degradation and T7 polymerase-assisted unbiased amplification. We demonstrate that TEAM-seq delivered, to the best of our knowledge, the highest reported coverage(70% for 100pg, 35% for 20pg) of single cell genomes in whole-genome 5mC sequencing.

**Introduction**

In recent years, the rapid development of single-cell sequencing technology has provided many valuable insights into complex biological systems, for example, revealing complex and rare cell populations and tracking the trajectories of distinct cell lineages in development [1]. Numerous single-cell RNA and DNA sequencing methods have been developed [2]. Heritable phenotype changes occur not only owing to changes in DNA and RNA nucleotide sequence, but also via epigenetic modifications, which do not change the DNA sequence itself. Among the epigenetic modifications, methylation of cytosine with the formation of 5-methylcytosine (5mC) is the most abundant epigenetic base change in vertebrates. This epigenetic modification has been intensively studied in relation to embryonic development, genomic imprinting, X-inactivation, cellular proliferation, and differentiation [3, 4]. Conventional approaches

1    to studying 5mC epigenetic marks usually rely on chemical treatment to convert

2    modified/unmodified sites to readable mutated nucleotides, which requires high DNA input [5].

3    However, the limited number of DNA copies significantly complicates the application of these

4    conventional approaches to studies in single cells.

5    The standard and well-established technology to detect sequence-specific 5mC marks is

6    bisulfite sequencing [6, 7], which utilizes sodium bisulfite to convert the unmethylated cytosine to

7    uracil after constructing the sequencing library. Methylated cytosine is resistant to bisulfite

8    conversion. Although bisulfite sequencing has been established as the gold standard for bulk

9    DNA methylation analysis, single-cell adaptations of this method for low DNA input face a

10   major hurdle of bisulfite-induced DNA degradation. Pioneer studies of single-cell bisulfite

11   sequencing were performed by integrating all conventional steps, including bisulfite

12   sequencing, into a single-tube reaction, followed by PCR amplification and deep sequencing [6,

13   7]. The constructed DNA library is usually broken by the harsh bisulfite treatment, and the

14   CpG/genome coverage drops to 5%/4% [8]. The CpG coverage rate can be improved to

15   approximately 18% through post-bisulfite adaptor tagging [9]. Single-cell genome-wide bisulfite

16   sequencing, in which the adaptor tags bisulfite-treated DNA with a 3′ stretch of random

17   nucleotides was shown to further improve the coverage by up to 5–48.4% of all CpG sites [10-13].

18   Recent advances in single-nucleus methylcytosine sequencing have been based on this

19   approach, with an additional use of adaptase and random priming to tag bisulfite-treated DNA

20   with adaptors [14].

21   In addition, recent advances have facilitated high-throughput single-cell methylation

22   sequencing. Single-cell combinatorial indexing for methylation analysis uses transposons to

23   assign combinatorial indices on single cells and then tag the other adaptor by random priming

3

1   after bisulfite conversion [15]. Although these previous methods improved the adaptor tailing

2   efficiency, they suffered from DNA degradation because of the harsh bisulfite treatment.

3   Recently, a gentle conversion method that combines ten-eleven translocation (TET) oxidation

4   with pyridine borane sequencing was reported [16], which used a high concentration of TET1 to

5   convert 5mC to 5CaC, which in turn, was converted to uracil by bicarbonate treatment.

6   However, high concentrations of TET1 are not available in most laboratories. Overall, single-

7   cell methylation sequencing coverage was poor, as it covered only 20% of CpG sites and less

8   than 10% of the genome. Therefore, it is necessary to improve the coverage, accuracy, and

9   read length of single-cell methylation sequencing.

10  The methods mentioned above improve the adaptor tagging efficiency or minimize DNA loss,

11  but they ignore biased PCR amplification that causes uniform genome coverage. Here, we

12  describe the T7-assisted enzymatic methyl-sequencing method (TEAM-seq) that improves

13  coverage. In TEAM-seq, transposons tagmentate genomic DNA with T7 adaptors. Then, a

14  gentle and high conversion efficiency kit is used to protect the methylation sites by TET2 and

15  convert the unmethylated sites to uracil by the activity of apolipoprotein B mRNA editing

16  enzyme catalytic subunit 3A (APOEC3A). This method can also be easily scaled up to high-

17  throughput analysis of single cells. With TEAM-seq, we achieved a superior 70% coverage

18  efficiency with only 2× sequencing depth and ultra-low DNA input, which, to the best of our

19  knowledge, is the highest single-cell methylation coverage ever reported.

20

21  **Results**

22  **Feasibility of enzyme-based conversion to detect the methylome**

4

1    To detect the methylome, the conventional method is bisulfite sequencing (BS), which relies on

2    the bisulfite-induced conversion of cytidine to uridine in the genome. However, 5-

3    methylcytidine (5mC) does not get converted in this way (Sfig.1a). Furthermore, this harsh

4    treatment usually results in very short fragments and significant loss of genomic DNA [17]. To

5    avoid the harsh bisulfite treatment, we consider using the gentle enzyme-based conversion

6    methods, such as TET-assisted pyridine borane sequencing (TAPS)[16] or commercially

7    available EM-seq (NEBNext® Enzymatic Methyl-seq Kit | NEB). Unlike directly converting 5mC

8    to U by TAPS, the commercially available enzyme-based method termed Enzymatic Methyl-

9    seq (EM-seq) has been developed recently by New England Biolabs, which uses APOEC3A to

10    convert cytidine to uridine for a gentler treatment of DNA [18]. The conversion protocol included

11    two general steps: 1) TET2 transformed 5mC/5hmC to 5-formylcytosine (5fC), 5-

12    carboxycytosine (5caC), and 5-(β-glucosyloxymethyl)cytosine (5gmC), which could be

13    protected from downstream APOEC3A conversion; 2) APOEC3A converted all the remaining

14    unmodified cytidines to uridines (Sfig.1a). The unconverted "C" in sequencing data

15    represented methylation sites, as in BS. A comprehensive analysis of the catalytic activity of

16    this enzyme and a comparative study of the NA12878 dataset have been published

17    previously[19]. Here, we sought to further extend the application of EM-seq to plant genomic

18    DNA of *Arabidopsis*.

19    The BS and EM-seq libraries were prepared in a pairwise manner. We used 1μg of

20    *Arabidopsis* DNA for the initial test. Both BS and EM-seq require PhiXDNA to balance the

21    sequencing signal bias. The sequencing quality and data yield were comparable between EM-

22    seq and BS (Sfig1.b–f). EM-seq was characterized by a **14%** lower duplication rate (Sfig.1c)

23    and better methylome coverage (Sfig.1d, e) compared to those observed with the BS

1   approach. To assess the accuracy of EM-seq, we added two different types of DNA to

2   estimate the false-positive rate (lambda DNA) and 600 bp synthetic methylated DNA to

3   estimate the false-negative rate. We found that false positive and false negative rates were

4   comparable for EM-seq and BS approaches (1.69% vs. 0.96% and 6.03% vs. 2.93%,

5   respectively) (Sfig.1c). The conversion rate of unmodified C in EM-seq was approximately

6   **98.31%** (Sfig.1c). The same software and pipeline were used to analyze the outputs of BS and

7   EM-seq experiments (Supplementalmaterial1_detailedprotocol).

8   Next, we analyzed the consistency of EM-seq and BS methylome data. At least four reads

9   covered more than 15 million CpG positions in both methods. We defined a base as

10   methylated if its methylation ratio was greater than 0.5 and found that 98.5% of CpG regions

11   (covered > 4 times) had a consistent modification pattern in EM-seq and BS experiments. (Fig.

12   1a). For example, both EM-seq and BS reported highly correlated methylation sites across

13   chromosomal regions (Fig. 1b). By further comparing the modification level for each 5mC

14   covered by at least **three** reads, we observed a good correlation (0.89) between EM-seq and

15   BS data (Fig. 1c). The density of the methylation across chromosome 4 or on gene elements

16   were similar between the two methods (Fig. 1d–f). For example, both methods showed equal

17   methylation level distribution around the CpG islands (CGIs) (Fig. 1g). Together, these results

18   indicate that EM-seq can directly replace whole genome BS (WGBS) and provide comparable

19   results. We then further examined whether this non-destructive method could be used for

20   methylation detection in ultra-low input samples, using various conversion methods.

21

22   **DNA polymerase amplification of a large uridine-rich fragment**

1   Multiple displacement amplification (MDA) is one of the most popular choices for non-

2   destructive amplification of the genome in ultra-low input samples [20, 21]. We started by

3   performing EM-MDA-seq in a 1-ng sample of genomic DNA (Fig. 2b). Genome cytidines were

4   converted to uridines by EM-seq, producing long uridine-rich genomic fragments (Fig. 2a). The

5   converted DNA fragments underwent MDA and 10-fold amplification was achieved after 30 min

6   of incubation (Fig. 2b). Subsequently, conventional DNA library construction was performed,

7   and EM-MDA-seq of this 1-ng input showed excellent performance (**84.79%** coverage at the

8   **6.96×** depth) (Fig. 2c, green). However, in the experiment with a picogram input, EM-MDA-seq

9   showed a large amplification bias and insufficient coverage (**5.41%** coverage at the **2.62×**

10  depth) (Fig. 2c, green). We then decided to lower the EM-MDA-seq bias with picogram input

11  samples by improving the priming efficiency and polymerase speed.

12  It has been previously reported that random hexamer priming results in a bias in nucleotide

13  composition and that this bias influences the uniformity of the location of the reads. The novel

14  MDA "Trueprime" method, which is based on the use of DNA primase without a random

15  hexamer, demonstrated superior breadth and uniformity of genome coverage with high

16  reproducibility [22] (Fig. 2b). However, Trueprime was not very efficient in amplifying the uridine-

17  rich converted DNA and yielded only **2.03%** coverage at the **3.94×** depth.  (Fig. 2c, blue).

18  It has been noted previously that DNA polymerase has base-pairing preference for nucleotide

19  incorporation [23], which may cause amplification bias during EM-MDA-seq. Most DNA

20  polymerases incorporate GT and AC nucleotides with high and low efficiency, respectively [23].

21  In our experiments, we also found that the speed of Phi29 was 5-fold lower in bisulfite-

22  converted DNA than in untreated genomic DNA (based on BGI production center reports),

23  supporting the notion of the low efficiency of adenosine incorporation. We tried several

1  solutions to normalize the nucleotide incorporation speed, including poly (ethylene

2  glycol)/trimethylammonium chloride solution, slow down amplification, and minimizing the bias.

3  We also tried an alternative robust polymerase, BST. However, none of these approaches

4  showed superiority in terms of uniform coverage (Fig. 2c). However, to solve the issue of the

5  AU bias in EM-seq, we used the GpC methyltransferase M.CviPI to pretreat the genome,

6  which preserves 25% more GC content through the EM-seq conversion (Fig. 2a). Surprisingly,

7  GpC methyltransferase treatment did not lead to obvious advantages over the control group in

8  a side-to-side comparison (Fig. 2c orange).

9  Thus, the DNA polymerases Phi29 and BST had amplification bias in the converted uridine-

10  rich DNA and were concluded to be unsuitable for EM-seq amplification. We then pursued an

11  alternative amplification method with a more even coverage, by using T7 RNA polymerase

12  amplification.

13

14  **Uniform amplification of ultra-low input uridine-rich DNA fragments by T7 RNA**

15  **polymerase**

16  Considering that the AC/GT bias amplification may be general problem for DNA polymerase,

17  we used RNA polymerase instead of DNA polymerase to avoid the bias amplification. An

18  improved single-cell whole genome amplification method based on the use of T7 polymerase

19  has been recently reported that outperformed other MDA approaches, enabling a more uniform

20  coverage of the genome [24]. Unlike the MDA, using Phi29 DNA polymerase to ceaselessly

21  synthesize the complementary strand on the single strand DNAs, T7 polymerase generated

22  marvelous RNAs on the DNA templates with less bias. Inspired by that finding, we have

23  developed a novel method, which we termed T7-Assisted Enzymatic methyl-seq (TEAM-seq).

1    In this method, we used Tn5 transposons to tagmentate the genome with biotinide adaptors

2    and obtain medium-sized fragments (500–2,000 bp). The tagmentated DNA was then

3    enzymatically converted, and single stranded DNAs were extended on the T7 promoter oligo to

4    generate the T7 promoter on the DNA (Fig. 3a). Throughout the process, we used streptavidin

5    beads to bind DNA without elution to minimize DNA loss in the purification. Finally, the

6    enzymatically converted DNA was amplified using T7 RNA polymerase. The final constructed

7    library was sequenced by the nanopore or Pacbio method to generate long reads for

8    paternal/maternal phasing (the detailed protocol and experimental information can be found in

9    the supplementarymaterial1). We used 100-pg or 20-pg of genomic DNA for the initial analysis.

10   The coverage and other parameters were better or comparable to those achieved by the

11   WGBS of a **1-μg** sample (**11.58%** duplication rate, coverage at **2.58×** depth). The consistency

12   between TEAM-seq and WGBS approaches was approximately ~93% (Fig. 3b). We also

13   summarized the features of the different published methods (Fig.3c). We then performed a

14   detailed study of the coverage and accuracy of TEAM-seq (comparing with scWGBS, and

15   conventional WGBS).

16   The sequencing quality and data yield for 100-pg/20-pg samples were similar to those

17   achieved by the conventional nanopore DNA sequencing. Analysis of the 20-pg input showed

18   a higher base volatility than that of the larger, 100-pg input (Sfig. 2a). The 20-pg library

19   resulted in **37.09%** genome coverage with **2.05×** sequencing depth, which are, to the best of

20   our knowledge, the highest recorded values for such analyses (Fig. 3b). The accumulation

21   curve showed the coverage saturation in each method (Sfig. 2b). For the picogram-level input,

22   TEAM-seq showed lower GC bias and equal coverage of CGIs(CpG islands) than those

23   achieved by the analyses of a 1-μg input by WGBS and 100-ng input by single cell WGBS

1 (Sfig. 2c–e). The higher and more uniform genome coverage resulted in a larger number of all

2 C sites covered by TEAM-seq in the case of a picogram input. In the experiment with a 20-pg

3 input, TEAM-seq covered 26.97% of CpG and 24.75% of C sites at the **2×** depth, whereas for

4 a 100 pg input, these values were higher: 54.85% and 52.79%, respectively. Further, TEAM-

5 seq covered 34.12% and 63.74% of CGIs in experiments with 20 and 100-pg inputs,

6 respectively.

7 TEAM-seq also showed a high consistency and methylation level correlation with the results of

8 the gold standard WGBS. Over ~**93%** of CpGs in TEAM-seq (100-pg/20-pg) showed

9 modification states that matched those detected by WGBS (Sfig. 3a). The correlations between

10 WGBS and TEAM-seq results (**100-pg**/**20-pg**) were **0.85/0.81** for the sites with >**20×** coverage.

11 Distributions of the methylated sites were also very similar (average standard error ± 0.2 on

12 CpG; ± 0.6 on CHG) on chromosome and gene element levels and close to those observed for

13 WGBS (Sfig. 3c, d). The methylation ratios of the 20 pg genomic DNA sample or single-cell

14 genomic DNA, which only had 1 or 2 genome copies, showed that the methylation distribution

15 was polarized (peaks on 0 or 1) (Sfig. 3b). For example, in the 20 pg input sample, a similar

16 methylation level distribution in the CGI was observed (Sfig. 3e). Moreover, we found that

17 TEAM-seq had a shallow coverage-related methylation ratio bias (Sfig. 3f). Three technicians

18 performed the experiments, and the reproducibility was sufficiently good (Sfig.4).

19 In summary, TEAM-seq uniformly amplified enzymatically converted DNA, resulting in better

20 coverage of samples with ultra-low input DNA amount. Next, we used TEAM-seq to perform

21 single-cell methylome analysis to expand its application scope.

22

23 **Indexed TEAM-seq for mid to high throughput single-cell analysis**

1   Using transposons, we designed a unique, relatively long barcode, which could be identified by

2   nanopore sequencing (Sfig. 5). After cell nuclei were sorted into plates, we labeled the cell

3   genome using the barcoded transposons. We then performed single-cell TEAM-seq on the

4   pooled cells (Please see Supplemental material 1_detailed protocol). Due to the throughput

5   limitation of nanopore sequencing, we only processed 50–100 cells. Our proposed method can

6   be easily scaled up and down based on your nanopore sequencing budget. We observed a

7   meager collision rate of **5%**, based on the artificial mixture of the human(HeLa) and

8   mouse(4T1) cells (Fig. 4a, b).

9   Next, we profiled an artificial mixture of HeLa cells and the HEK293T cells. We could

10  debarcode >70% reads and **64 cells** (100% cell capturing efficiency), all of which passed the

11  quality filters (depth > 0.2×, mapping rate > **45%**, genome coverage > **3%**). Each cell had a

12  median of 10 million reads and a mean alignment rate of **56%**, approaching the levels seen in

13  the TEAM-seq experiment with 20-pg DNA input (Sfig. 6a). These data translate to the

14  coverage of mappable CpG ranging from ~5% to ~20% (Sfig. 6c) and genome coverage

15  ranging from ~10% to ~30% (Sfig. 6a). As expected, the sequencing depth was not yet

16  saturated, and an increase in the amount of data could further improve the genome coverage

17  (Sfig. 6b). Among these covered CpG and genome sites, ~0.2 million Ensembl builds were

18  covered by at least half of the cell population (Sfig. 6f). We did not observe any abnormal

19  coverage bias on the whole genome scale or in CGI (Sfig. 6e, d). Next, we summarized the

20  methylation status of each cell across Ensembl builds to observe the accuracy and variance of

21  methylation mark detection. The correlations between single cells ranged from **0.4 to 0.7** in

22  0.1–2 kb fragments of Ensembl build (Sfig. 7a). The variance across the genome comprised

23  **0.01–0.05 in CGIs, enhancers, and gene bodies** (Sfig. 7b). On

1    H3K27me3/H3K27ac/H3K4me3, all cells produced the expected nucleosome pattern with

2    corresponding heterochromatin/euchromatin properties (Sfig. 7c).

3    Next, we profiled a mixture of HEK293T and HeLa human cell lines (Fig. 4). We performed

4    non-negative matrix factorization (NMF) followed by k-means embedding to project cells in

5    two-dimensional space, producing two clearly defined clusters (Fig. 4c, d). We correlated the

6    methylation rate of the collapsed cluster with publicly available WGBS datasets (Fig. 4e). One

7    cell cluster showed the highest correlation with the corresponding cell type (Fig. 4e). To display

8    the advantage of increased coverage of the single cell, we down-sampled the coverage and

9    found that the identified two cluster distances decreased with coverage until merging (Fig. 4d).

10    Overall, TEAM-seq could generate high methylation site coverage on the indexed single cells

11    in a high-throughput manner.

12

### Discussion

14    In the current study, we developed the TEAM-seq method based on the enzymatic cytidine

15    conversion (EM-seq) and T7 polymerase amplification, achieving a very high genome

16    coverage in 100 pg and 20 pg input samples, and even in single cells. However, EM-seq still

17    demonstrated 0.7% lower conversion rates than the gold standard BS-seq, resulting in a

18    higher false-positive/negative rate. The higher number of false positive/negative methylation

19    marks may be caused by the substrate/motif preference of Tet1/2 and APOEC3A or different

20    enzymatic efficiency. The slightly higher false positive/negative rate would be problematic in

21    single-cell analysis. Because a single cell has only one/two copies of genome, the accuracy is

22    difficult to correct. However, this problem cannot be easily avoided in enzymatic conversion.

23    Given the significant DNA loss during the harsh bisulfite treatment, we believe that enzymatic

1    conversion is a better choice for single-cell methylation analysis. Moreover, we are developing

2    a machine learning algorithm that may help distinguish false positive/negative methylation

3    marks in single cells.

4    We have tried many other approaches in the current study, but most methods produced

5    surprisingly unsatisfactory results. The Trueprime kit demonstrated superior performance in

6    single-cell genome amplification compared to that achieved with MDA, but not for

7    enzymatically converted uridine-rich DNA. We suspected that heavy modifications

8    (5CaC/5fC/5gmC) generated by EM-seq inhibited primer synthesis by primase. A similar

9    phenomenon was also observed during the GpC methyltransferase treatment, which

10   preserved 25% of the GC content in the enzymatic conversion. The high GC content adversely

11   affected MDA owing to the parallel increase in heavy modifications (5CaC/5fC/5gmC). We also

12   observed a related phenomenon in a mixture of HEK293T and HeLa cells. HEK293T cells

13   generated twice more reads than HeLa cells, suggesting the preferential amplification of

14   HEK293T cell DNA. This bias was not observed in a mixture of HeLa and 4T1 cells.

15   Considering the higher abundance of 5mC marks in cancer cells[25], heavy modifications may

16   also lead to the inefficiency of T7 polymerase. Therefore, the vast global difference in the

17   methylation rate may cause biased amplification in cell populations for pooling the indexed

18   cells.

19   In TEAM-seq, we used Tn5 to fragment genomic DNA, but the fragment length was difficult to

20   control precisely in single cells. The varied fragment length required nanopore or Pacbio

21   sequencing. The nanopore method has higher error rate than Pacbio. However, we found that

22   two reads per site could confidently correct the sequencing errors. Otherwise, our patented

1  linked-reads method [26] on the Pacbio platform could be an alternative to nanopore sequencing,

2  which could increase 5x output of Pacbio sequencing.

3  Overall, TEAM-seq reported genome coverage for ultralow input DNA amounts and DNA from

4  single cells. This method could help in studies of samples with a limited DNA amount, e.g.,

5  embryos, oocytes, clinical tissues, etc. The high coverage of single cell genomes may further

6  disclose methylation heterogeneity of different cell populations.

7

8  **Methods**

9  **Experimental methods**

10 **Assess the EM-seq and WGBS**

11 *Arabidopsis* DNAs were given by the other lab. The EM-seq was performed and NEB EM-seq

12 protocol (NEB E7120S)(1ug DNA per reaction). The WGBS was parallelly performed by using

13 MGIEasy Whole Genome Bisulfite Sequencing Library Prep Kit (MGI 1000005251). Both

14 methods were similar except the conversion steps. The final products were sequenced by

15 MGISEQ-2000.

16 **MDA for the converted DNA**

17 Both of the WGBS and EM-seq converted cytidine for the constructed DNA library, which was

18 fragmented and ligated with adaptors. Unlike the WGBS, MDA-EM-seq used the Phi29 to do

19 the multiple displacement amplification on the non-fragmented DNA. The genomic DNA were

20 directly processed to EM-seq conversion steps

21 (https://international.neb.com/protocols/2019/03/28/protocol-for-use-with-standard-insert-

22 libraries-370-420-bp-e7120, 1.5-1.9). The converted single-strand DNAs were around 6~10kb

23 after enzyme treatment. The clean-up single-strand DNAs were processed to the MDA

1  amplification by REPLI-g Single Cell Kit (Qiagen 150343). For Trueprime MDA-EM-seq, the

2  clean-up DNAs underwent the protocol SYGNIS TruePrime™ WGA & Single Cell WGA Kits

3  (Lucigen), instead of the REPLI-g single cell kit. In the EM-MDA-seq with additives, all the

4  addictive (3% PEG,100mM TMAC, 3%PEG+100mM TMAC) were added to the REPLI-g Single

5  Cell reaction. In the EM-MDA-seq with BST amplification, we used the BST2.0 (NEB M0537S)

6  to substitute the Phi29 in REPLI-g single cell kit. The amplified DNAs were processed as

7  MGIEasy Universal DNA Library Prep Kit V1.0 and sequenced in MGISEQ-2000.

8  **scWGBS_ng**

9  The scWGBS_ng was downloaded [27]. Briefly, bisulfite conversion was performed using the EZ

10  DNA Methylation-Direct Kit (Zymo Research, D5020) according to the manufacturer's protocol,

11  with the modification of eluting the DNA and minimizing the DNA loss. The single stranded library

12  construction was followed.

13  **TEAM-seq**

14  We first fragment the DNA by Tn5 transposon with biotin adaptor. After tagmentation, the DNA

15  with biotin adaptors could be bound with the streptavidin beads. Then the DNAs carried on

16  streptavidin       beads       were       processed       to       EM-seq

17  (https://international.neb.com/protocols/2019/03/28/protocol-for-use-with-standard-insert-

18  libraries-370-420-bp-e7120, 1.5-1.9). The DNA immobilized on the beads could be smoothly

19  transfer to the next reaction without the significant loss. After conversion, the single strand DNAs

20  (beads-on reaction) were extended on the complementary oligos to generate the T7 promoters.

21  Then we performed the bead-on T7 transcription reaction by HiScribe™ T7 High Yield RNA

22  Synthesis Kit (NEB E2040S). The amplified RNAs were processed to the PCR-cDNA

23  sequencing kit (nanopore). The detailed description of methods can be found in Supplemental

1    Material supplementalmaterial1_detailedprotocol. For the single cell TEAM-seq, we used the

2    Qiagen protease to lysis the single cell to release DNAs and heat inactivated the Qiagen

3    protease before processing to the Tn5 tagmentation.

4    **Cell culture**

5    Human mammary gland carcinoma cell line HeLa/HEK293T were obtained from ATCC. MCF-7

6    were grown in DMEM (Gibco,11995065) supplemented with 10% FBS (Gibco,10099141),

7    0.01mg/ml insulin (HY-P1156, MedChemExpress), and 1% penicillin-streptomysin (Gibco,

8    15140122). Cell line was regularly checked for mycoplasma infection (Yeasen, 40612ES25).

9    **Bioinformatics method**

10    **Data processing**

11    Sequencing adapters were cut using fastp (v0.21.0). The remaining reads were trimmed based

12    on their theoretical barcode locations. Then they were aligned to the barcode white list in the

13    following pattern: Ns-GGGAGATACAACCTACAATCACT-10bp barcode-

14    AAATATATATAAAAAACAA-Ns using bowtie2 local alignment mode (v2.3.4). Only reads

15    aligned to the barcode reference with high mapping quality would be processed and used for

16    subsequent analysis.

17    To assess the accuracy of demultiplexing, 10,048 reads were simulated with a length of 52 bases,

18    including a 23bp 5' primer, 10bp barcode from the white list, and a 19bp 3' primer. In order to

19    intimate nanopore sequencing errors, a mean indel rate of 2.96% and a mean mismatch rate of

20    2.77% were introduced to the simulated data.

21    Reads were aligned by methylpy (v1.3) in the following steps: (1) reads mapped to the forward

22    strand GRCh37 reference genome were de-duplicated with the parameter "--pbat --remove-

23    clonal --path-to --picard --min-mapq 15". (2) Remaining unmapped reads were then mapped to

the reverse strand of the reference genome in the same way. (3) Merge all mapped bam, then call methylation status by methylpy call_methylated_sites. In terms of human-mouse mixed library, a combined human-mouse hybrid genome was used instead of GRCh37 reference.

Fastq files of scWGBS,sci-MET, snmC-seq and scBS-seq downloaded from the database were aligned to the reference genome and called the methylation ratio using Bismark (v0.22.3). Any other conventional WGBS data was analysed using our standard pipeline.

**Analysis of cytosine modifications called by TEAM-seq, scWGBS and WGBS**

Chromosomal methylation status was plotted using CG/CHG/CHH positions covered by more than three times. The chromosome was then binned into equal sizes (100kb for Arabidopsis Chr4, 2Mb for human Chr1 and Chr4). Mean methylation ratio per bin was calculated and plotted along the chromosome. The mean occurrences of C/CG/CHG/CHH in each bin was computed and plotted in the same way.

Information of histone (H3K27me3/H3K27ac/H3K4me3) ChIP-seq peaks was downloaded from the ENCODE database. Bed files of gene elements including CGIs, TSS, TES and gene body were downloaded from UCSC genome browser. These regions were binned into 20 or 50 windows, flanking by 50 windows of the up-/downstream regions. The methylation ratio of each bin was averaged and plotted using R.

**Average coverage depth in CGI**

Total nucleotide count for each bin was reported using samtools bedcov (v1.9), then average coverage depth in each window of the CGI region was computed for TEAM-seq, scWGBS and WGBS. To overcome differences in sequencing depth, TEAM-seq and WGBS data were down-sampled to the corresponding sequencing depth of scWGBS. Average coverage depth for single cell data was computed in the same manner.

**Correlation and consistency analysis**

For each single cell, weighted methylation ratios of the Ensembl Regulatory Build regions were computed in a previously described method [28] [10]. We then calculated the Pearson correlation of regions including all Ensembl Regulatory loci, promoter region, enhancer region and CTCF binding sites. Clustered heatmap was then plotted using pheatmap in R.

In order to correlate with bulk WGBS data, we merged the methylation information of all cells in cluster 1 and calculated methylation ratio at single nucleotide level. The methylation ratio was then correlated to WGBS datasets available from the ENCODE database as well as the conventional HeLa WGBS data [29].

The consistency of methylation status was compared between WGBS and EM-seq, TEAM-seq 100pg and 1ug WGBS, TEAM-seq 20pg and WGBS, and among three TEAM-seq 20pg libraries. For each comparison, only positions covered more than three times were used and a methylation ratio of 0.5 was used as the threshold. For example, if the methylation ratio of a cytosine is either below 0.5 or above 0.5 in both libraries, the methylation pattern was considered as consistent. Venn diagram and Euler diagram were plotted using VennDiagram and eulerr packages in R.

**NMF decomposition and k-means clustering**

Non-negative Matrix Factorization (NMF) is a decomposition algorithm to split a matrix V into two matrices W (dictionary matrix) and H (meta-feature matrix). Due to its clustering property, it has been widely used in DNA methylation studies to identify features that contribute the most to the cluster. Therefore, using Ensembl Regulatory Build as the feature, NMF clustering was performed to distinguish HeLa and HEK293T cell lines. Then, k-means clustering algorithm was used to cluster cells by selecting 100%, 60%, 30% and 10% features to investigate the trend of two clusters getting more similar to each other. All clustering analyses were performed in R.

18

1

2

3

**Data availability**

TEAM-seq_100pg, TEAM-seq_20pg, HeLa&4T1 mixture data as well as HeLa&HEK293T mixture data are available at China National GeneBank (CNGB) [30, 31] with project number of CNP0002270.

**Author contributions**

CT designed and supervised the experiments. JW performed the experiments; YTF, FC, and CT performed bioinformatics data analysis. All authors collectively analyzed experimental data. All authors read and approved the final draft of the manuscript.
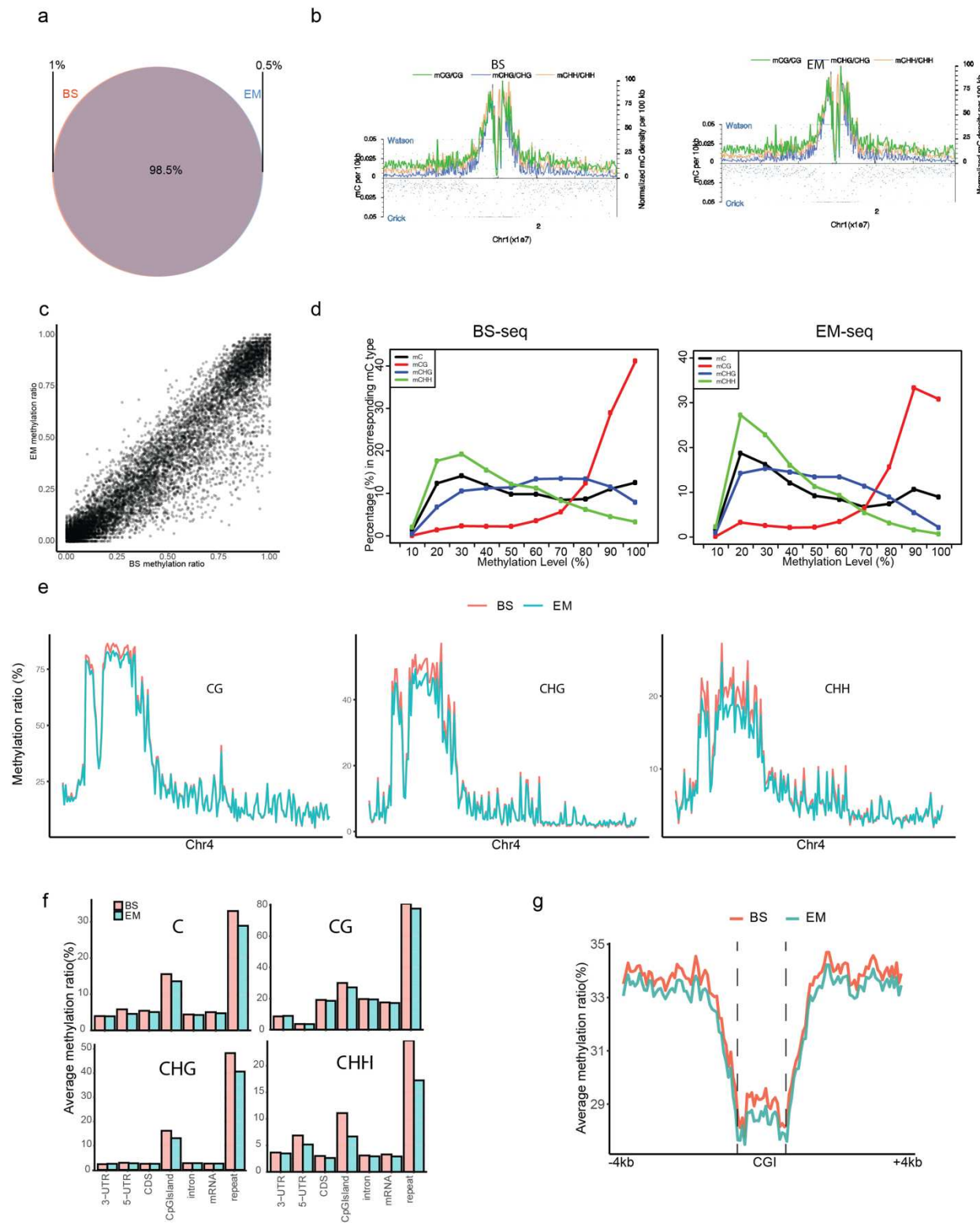
**Competing interest**

A patent application has been filed by BGI Genomics Co Ltd for the technology disclosed in this publication.

# Bibliography

1.  Hwang, B., Lee, J.H. & Bang, D. Single-cell RNA sequencing technologies and bioinformaticspipelines. *Experimental & Molecular Medicine* **50**, 1-14 (2018).
2.  Nawy, T. Single-cell sequencing. *Nature Methods* **11**, 18-18 (2014).
3.  Li, E. Chromatin modification and epigenetic reprogramming in mammalian development. *Nat Rev Genet* **3**, 662-673 (2002).
4.  Bird, A. Perceptions of epigenetics. *Nature* **447**, 396-398 (2007).
5.  Tollefsbol, T.O. Methods of epigenetic analysis. *Methods Mol Biol* **287**, 1-8 (2004).
6.  Guo, H. et al. Profiling DNA methylome landscapes of mammalian cells with single-cell reduced-representation bisulfite sequencing. *Nat Protoc* **10**, 645-659 (2015).
7.  Guo, H. et al. Single-cell methylome landscapes of mouse embryonic stem cells and early embryos analyzed using reduced representation bisulfite sequencing. *Genome Res* **23**, 2126-2135 (2013).
8.  Karemaker, I.D. & Vermeulen, M. Single-Cell DNA Methylation Profiling: Technologies and Biological Applications. *Trends in Biotechnology* **36**, 952-965 (2018).
9.  Miura, F., Enomoto, Y., Dairiki, R. & Ito, T. Amplification-free whole-genome bisulfite sequencing by post-bisulfite adaptor tagging. *Nucleic Acids Research* **40**, e136-e136 (2012).
10. Smallwood, S.A. et al. Single-cell genome-wide bisulfite sequencing for assessing epigenetic heterogeneity. *Nature Methods* **11**, 817-820 (2014).
11. Miura, F. & Ito, T. Highly sensitive targeted methylome sequencing by post-bisulfite adaptor tagging. *DNA Research* **22**, 13-18 (2015).
12. Kobayashi, H. et al. Repetitive DNA methylome analysis by small-scale and single-cell shotgun bisulfite sequencing. *Genes to Cells* **21**, 1209-1222 (2016).
13. Farlik, M. et al. Single-Cell DNA Methylome Sequencing and Bioinformatic Inference of Epigenomic Cell-State Dynamics. *Cell Reports* **10**, 1386-1397 (2015).
14. Luo, C. et al. Single-cell methylomes identify neuronal subtypes and regulatory elements in mammalian cortex. *Science* **357**, 600-604 (2017).
15. Mulqueen, R.M. et al. Highly scalable generation of DNA methylation profiles in single cells. *Nat Biotechnol* **36**, 428-431 (2018).
16. Liu, Y. et al. Bisulfite-free direct detection of 5-methylcytosine and 5-hydroxymethylcytosine at base resolution. *Nature Biotechnology* **37**, 424-429 (2019).
17. Frommer, M. et al. A genomic sequencing protocol that yields a positive display of 5-methylcytosine residues in individual DNA strands. *Proc Natl Acad Sci U S A* **89**, 1827-1831 (1992).
18. Vaisvila, R. et al. EM-seq: Detection of DNA Methylation at Single Base Resolution from Picograms of DNA. *bioRxiv*, 2019.2012.2020.884692 (2020).
19. Vaisvila, R. et al. Enzymatic methyl sequencing detects DNA methylation at single-base resolution from picograms of DNA. *Genome Res* **31**, 1280-1289 (2021).
20. Dean, F.B., Nelson, J.R., Giesler, T.L. & Lasken, R.S. Rapid amplification of plasmid and phage DNA using Phi 29 DNA polymerase and multiply-primed rolling circle amplification. *Genome Res* **11**, 1095-1099 (2001).
21. Huang, L., Ma, F., Chapman, A., Lu, S. & Xie, X.S. Single-Cell Whole-Genome Amplification and Sequencing: Methodology and Applications. *Annu Rev Genomics Hum Genet* **16**, 79-102 (2015).

22. Picher, Á.J. et al. TruePrime is a novel method for whole-genome amplification from single cells based on TthPrimPol. *Nature Communications* **7**, 13296 (2016).

23. Choi, J.Y., Lim, S., Eoff, R.L. & Guengerich, F.P. Kinetic analysis of base-pairing preference for nucleotide incorporation opposite template pyrimidines by human DNA polymerase iota. *J Mol Biol* **389**, 264-274 (2009).

24. Chen, C. et al. Single-cell whole-genome analyses by Linear Amplification via Transposon Insertion (LIANTI). *Science* **356**, 189-194 (2017).

25. Ehrlich, M. DNA methylation in cancer: too much, but also too little. *Oncogene* **21**, 5400-5413 (2002).

26. Zheng, Y.-F. et al. HIT-scISOseq: High-throughput and High-accuracy Single-cell Full-length Isoform Sequencing for Corneal Epithelium. *bioRxiv*, 2020.2007.2027.222349 (2020).

27. Farlik, M. et al. Single-cell DNA methylome sequencing and bioinformatic inference of epigenomic cell-state dynamics. *Cell Rep* **10**, 1386-1397 (2015).

28. Zerbino, D.R., Wilder, S.P., Johnson, N., Juettemann, T. & Flicek, P.R. The ensembl regulatory build. *Genome Biol* **16**, 56 (2015).

29. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).

30. Chen, F.Z. et al. CNGBdb: China National GeneBank DataBase. *Yi Chuan* **42**, 799-809 (2020).

31. Guo, X. et al. CNSA: a data repository for archiving omics data. *Database (Oxford)* **2020** (2020).

1

1   Fig1. The data consistency between BS-seq and EM-seq. We annotated the confident

2   methylation sites with at least four reads covered and the methylation ratio is larger than 0.5.

3   98.5% methylated sites in BS-seq were detected in both EM-seq and BS-seq (a). We used the

4   Chr1 as an example, and we found that the methylated CG/CHH/CHG sites identified by these

5   two methods were distributed similarly across the Chr1 (b). We further looked at the

6   methylation level on each site. The global methylation ratio correlation is around 0.89

7   (spearman, $p<0.05$) (c). The methylation ratio density plot also showed the similar pattern

8   between two methods. The medium methylation level of CHH/CHG in EM-seq was slightly

9   lower than in the BS-seq (d). For example, the CG/CHG/CHH methylation ratio distribution

10  across Chr4 were similar between EM-seq and BS-seq. The CHH/CHG methylation level in

11  EM-seq demonstrated slightly lower level than in BS-seq in highly methylated regions (5%

12  lower) (e). By further looking at the gene elements, the average methylation ratios on the gene

13  elements were also slightly lower in EM-seq than in BS-seq, typically for CHH (f). More

14  detailed investigations on methylation patterns around CpG islands (CGIs) showed the

15  methylation pattern distributed similarly by using these two methods (g).

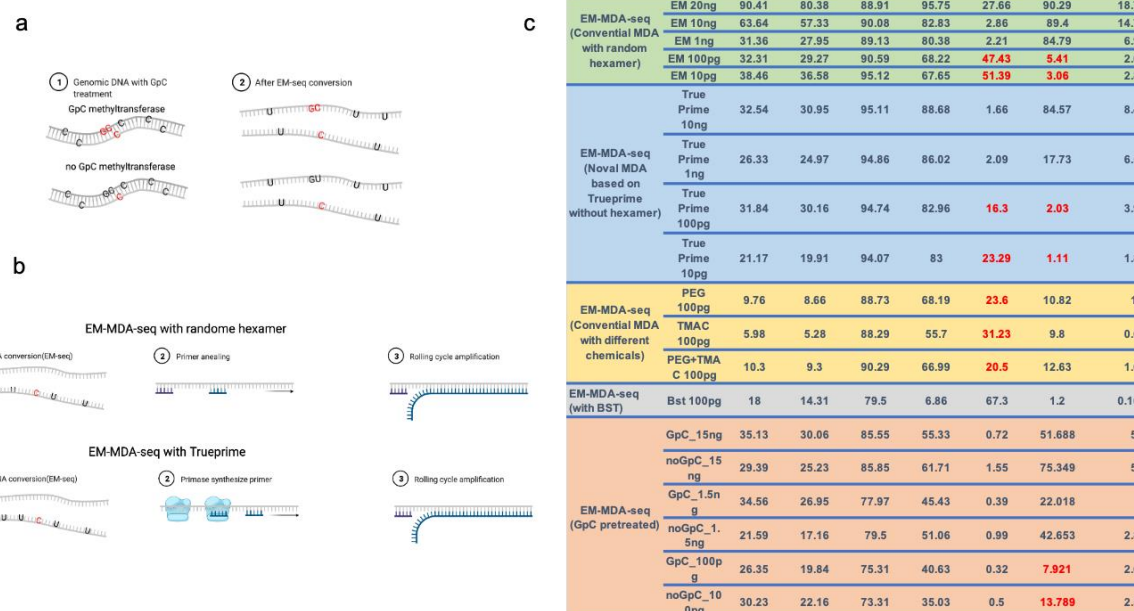| Method | Sample | Raw bases (G) | Clean bases (G) | Clean_rate (%) | Mapping rate (%) | Duplication rate (%) | Coverage (%) | Depth(X) |
|---|---|---|---|---|---|---|---|---|
| EM-MDA-seq (Convential MDA with random hexamer) | EM 20ng | 90.41 | 80.38 | 88.91 | 95.75 | 27.66 | 90.29 | 18.73 |
| | EM 10ng | 63.64 | 57.33 | 90.08 | 82.83 | 2.86 | 89.4 | 14.76 |
| | EM 1ng | 31.36 | 27.95 | 89.13 | 80.38 | 2.21 | 84.79 | 6.96 |
| | EM 100pg | 32.31 | 29.27 | 90.59 | 68.22 | 47.43 | 5.41 | 2.62 |
| | EM 10pg | 38.46 | 36.58 | 95.12 | 67.65 | 51.39 | 3.06 | 2.88 |
| EM-MDA-seq (Noval MDA based on Trueprime without hexamer) | True Prime 10ng | 32.54 | 30.95 | 95.11 | 88.68 | 1.66 | 84.57 | 8.48 |
| | True Prime 1ng | 26.33 | 24.97 | 94.86 | 86.02 | 2.09 | 17.73 | 6.56 |
| | True Prime 100pg | 31.84 | 30.16 | 94.74 | 82.96 | 16.3 | 2.03 | 3.94 |
| | True Prime 10pg | 21.17 | 19.91 | 94.07 | 83 | 23.29 | 1.11 | 1.88 |
| EM-MDA-seq (Convential MDA with different chemicals) | PEG 100pg | 9.76 | 8.66 | 88.73 | 68.19 | 23.6 | 10.82 | 1.5 |
| | TMAC 100pg | 5.98 | 5.28 | 88.29 | 55.7 | 31.23 | 9.8 | 0.67 |
| | PEG+TMAC 100pg | 10.3 | 9.3 | 90.29 | 66.99 | 20.5 | 12.63 | 1.65 |
| EM-MDA-seq (with BST) | Bst 100pg | 18 | 14.31 | 79.5 | 6.86 | 67.3 | 1.2 | 0.107 |
| EM-MDA-seq (GpC pretreated) | GpC_15ng | 35.13 | 30.06 | 85.55 | 55.33 | 0.72 | 51.688 | 5.5 |
| | noGpC_15ng | 29.39 | 25.23 | 85.85 | 61.71 | 1.55 | 75.349 | 5.1 |
| | GpC_1.5ng | 34.56 | 26.95 | 77.97 | 45.43 | 0.39 | 22.018 | 4 |
| | noGpC_1.5ng | 21.59 | 17.16 | 79.5 | 51.06 | 0.99 | 42.653 | 2.89 |
| | GpC_100pg | 26.35 | 19.84 | 75.31 | 40.63 | 0.32 | 7.921 | 2.67 |
| | noGpC_100pg | 30.23 | 22.16 | 73.31 | 35.03 | 0.5 | 13.789 | 2.57 |

Fig2. The MDA amplification on the EM-seq converted genome. After EM-seq conversion, most of the unmodified Cs were converted to U by APOEC3A deaminase. The GC content was decreased nearly 48%. The methylated Cs were protected by Tet2/enhancer, and stay as the cytidine analogs (5caC,5gmC,5fC) in the deamination process. The non-destructive method produced the pretty long uridine-rich ssDNAs (~2kb). We also used GpC methyltransferase to pretreat the genomic DNAs. Considering the GpCs were rare on mammalian genome, the artificial GpCs should not affect the native methylation identification (CpG,CHH,CHG). The methylated GpC could also be preserved in the deamination step (a). The MDA method used the random hexamers to prime on the uridine-rich ssDNAs from EM-seq. Phi29 extended on the hexamers and displaced the DNA ahead. Then the displaced ssDNA could be amplified again by random hexamers. The process was termed as "EM-MDA-seq" (b). The hexamers annealing efficiency is largely rely on the GC contents and annealing temperature. To avoid the random hexamer annealing bias, the novel method Trueprime, the

24

1    combination the TthPrimPol's unique ability to synthesize DNA primers with the highly

2    processive Phi29 DNA polymerase (Φ29DNApol) enables near-complete whole genome

3    amplification (b). The EM-MDA-seq of high input DNA (20~1ng) produce the normal genome

4    coverage around 80%. However, the EM-MDA-seq of picogram input DNA only give us 3~5%

5    genome coverage with very high duplication rate 47%. The Trueprime amplification also

6    resulted the similar low genome coverage 1~2% for picogram input DNA. The PEG have the

7    weakly positive effect to improve the genome coverage, and TMAC did not show this sign

8    either. We then tried use robust BST to replace Phi29 in MDA, which also produce the shallow

9    genome coverage 1.2%.  We then tried to use GpC methyltransferase to preserve more

10    cytidine in the deamination step by converting all the GpCs to artificially methylated GpCs(a).

11    By side-to-side comparison with the non-GpC treatment control, the GpC methyltransferase

12    treatment did not demonstrated the significant advantages in genome coverage from

13    nanogram input to picogram input. The methods were described in the
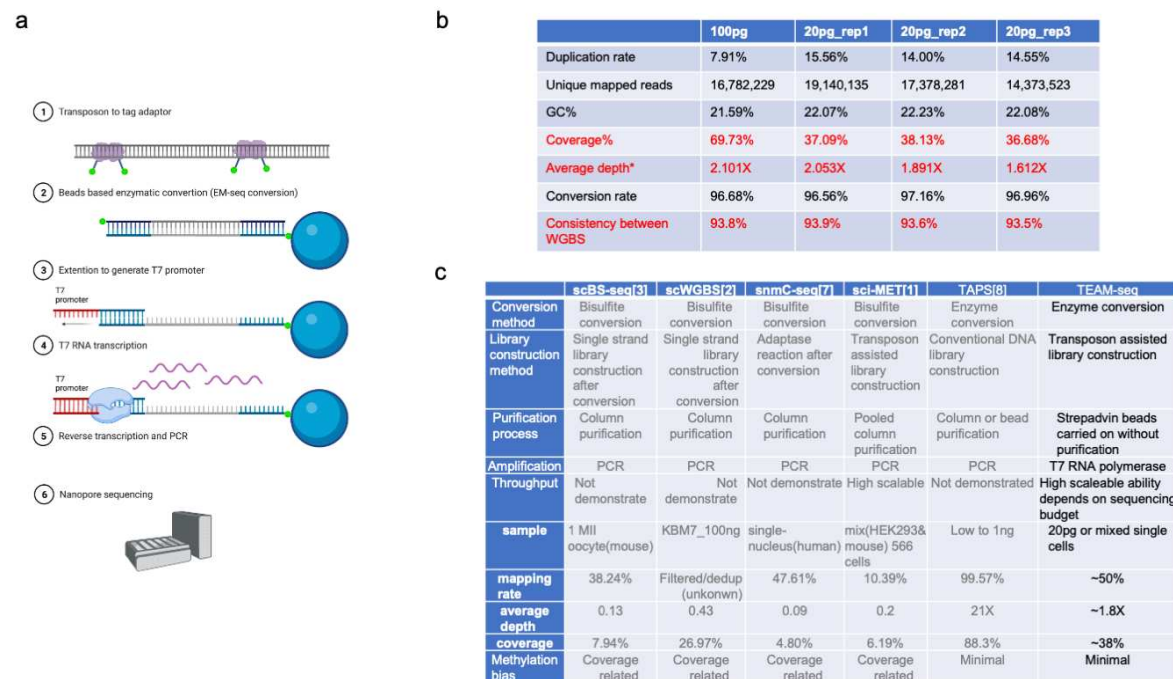
14    supplementalmaterial1_detailedprotocol part D.

**a**



**b**

| | 100pg | 20pg_rep1 | 20pg_rep2 | 20pg_rep3 |
|---|---|---|---|---|
| Duplication rate | 7.91% | 15.56% | 14.00% | 14.55% |
| Unique mapped reads | 16,782,229 | 19,140,135 | 17,378,281 | 14,373,523 |
| GC% | 21.59% | 22.07% | 22.23% | 22.08% |
| Coverage% | 69.73% | 37.09% | 38.13% | 36.68% |
| Average depth* | 2.101X | 2.053X | 1.891X | 1.612X |
| Conversion rate | 96.68% | 96.56% | 97.16% | 96.96% |
| Consistency between WGBS | 93.8% | 93.9% | 93.6% | 93.5% |

**c**

| | scBS-seq[3] | scWGBS[2] | snmC-seq[7] | sci-MET[1] | TAPS[8] | TEAM-seq |
|---|---|---|---|---|---|---|
| Conversion method | Bisulfite conversion | Bisulfite conversion | Bisulfite conversion | Bisulfite conversion | Enzyme conversion | Enzyme conversion |
| Library construction method | Single strand library construction after conversion | Single strand library construction after conversion | Adaptase reaction after conversion | Transposon assisted library construction | Conventional DNA library construction | Transposon assisted library construction |
| Purification process | Column purification | Column purification | Column purification | Pooled column purification | Column or bead purification | Strepadvin beads carried on without purification |
| Amplification | PCR | PCR | PCR | PCR | PCR | T7 RNA polymerase |
| Throughput | Not demonstrate | Not demonstrate | Not demonstrate | High scalable | Not demonstrated | High scaleable ability depends on sequencing budget |
| sample | 1 MII oocyte(mouse) | KBM7_100ng | single-nucleus(human) | mix(HEK293&mouse) 566 cells | Low to 1ng | 20pg or mixed single cells |
| mapping rate | 38.24% | Filtered/dedup (unkonwn) | 47.61% | 10.39% | 99.57% | ~50% |
| average depth | 0.13 | 0.43 | 0.09 | 0.2 | 21X | ~1.8X |
| coverage | 7.94% | 26.97% | 4.80% | 6.19% | 88.3% | ~38% |
| Methylation bias | Coverage related | Coverage related | Coverage related | Coverage related | Minimal | Minimal |

1

Fig3. The schematic procedure of the TEAM-seq and its overall performance. The transposons were assembled with Tn5 enzyme and the DNA adaptor, with biotin on the on the 5' end. After the tagmentation, the biotin adaptors were attached to the genomic DNAs and fragment it. After the denature the fragmented DNA were bound to the M280 strepadvin beads. Then the DNA carried on the beads went through the whole EM-seq protocol (detail could be found in method). Then the single stranded DNA were extended on the T7 oligos to generate the T7 promoters. After that, we performed the beads-on T7 transcription reaction for overnight. After the enough RNAs generated, the cDNA nanopore library were constructed and the detailed sequence information could be found on supplemental material_1. The final product could be sequenced on the nanopore(a). After we get the sequence, we analyze the data through the same protocol of the WGBS(1ug). The depth was calculated by average read count covered among all bases (ATGC). The whole genome coverage (the percentage of genomic sequence was covered by reads) is 69% for 100-pg TEAM-seq with 2x depth and 37% for 20-pg TEAM-

1    seq with 1.8x depth. We selected high confidence sites to calculate the consistency. The

2    methylated sites were defined as >0.5 methylation ratio with three covered reads, and the

3    unmethylated sites were defined as <0.5 methylation ratio. The ~93% genomic sites have the

4    matched methylation status in the WGBS, showing the good consistency between low input

5    TEAM-seq and the gold standard WGBS(b). (c) We summarized the genome coverages and

6    key features of the different methods from previous publications [1-4] [5] [6]. Column purification is

7    the key step to limit the scalability. scBS-seq, scWGBS and snmC-seq constructed the library

8    with barcodes after column purification. Therefore, the throughput of these three methods may

9    be limited. The sci-MET and TEAM-seq constructed the library with barcodes before

10   conversion, offering the simplicity in the downstream processing. The TEAM-seq took the

11   advantages of the transposon library construction, beads carried on reaction and T7 RNA

12   polymerase amplification to minimize the loss of the DNA and achieve uniform amplification.

13

14

Fig4. The processes and clustering results of the TEAM-seq on indexed single cells. The cells/nucleuses were sorted into the 96/384-well plate. Then the cells/nucleuses were lysed and tagmentate with indexed transposons. The barcoded genomic fragments were pooled and processed into TEAM-seq. The experimental process and bioinformatic process were described in Supplementalmaterial1_detailedprotocol(a). The barcode collision experiment

1   used 1:1 mixed 4T1/Hela cells and sampled 52 cells (corresponding 52 barcodes)(b). The two

2   barcodes, which did not pass the quality filter (Depth<0.2x), can be classified into neither

3   mouse nor human cells (red). Other barcodes could be classified into either mouse or human

4   cells (1:1 ratio), indicating the barcode collision rate <4%. The cell clustering experiment used

5   1:1 mixed HEK293T/Hela cells and sampled 52 cells (c,d). The results were summarized as

6   the matrix, including cells (rows) and corresponding methylation ratio in ensemble builds

7   (columns). By the algorithm Non-negative matrix factorization (NMF) [6], the matrix V is

8   factorized into (usually) two matrices W and H, with the property that all three matrices have no

9   negative elements (reduce dimension to 2 dimensions). We then studied the most suitable

10  cluster number (k value) for cell classification. As expected, the (k=2) two clusters are propriate

11  for classifying HEK293T and Hela cells. The increasing cluster number did not significantly

12  improve the classification(c). We then down sampled the data from 100% to 10%(d). The two

13  separated cell clusters gradually merged with the decreasing genome coverage, suggesting

14  the importance of the high genome coverage. Also the cell clusters showed the highest

15  correlation with the corresponding cell types from the database.

16

Sfig1. The quality control of the EM-seq and BS-seq. We showed the schematic plot of the

mechanism of EM-seq and BS-seq (a). Bisulfite treatment could convert all the unmethylated

cytidine(C) to uridine(U) without affecting methylated cytidine (5mC). The 5mC could be

preserved and read as "C" in the sequencing. EM-seq used the Tet2/enhancer to transform

the 5mC/5hmC to 5caC,5fC, and 5gmC, which could be preserved in the APOEC3A treatment.

1    APOEC3A then deaminate the cytidine to uridine, and cannot convert 5caC,5fC,5gmC to

2    uridine. The preserved "5mC" signal could also be recognized as "C" in sequencing. The

3    phenotype of these two types of sequencing data was similar, so the analysis pipelines are the

4    same. In the sequencing data, both EM-seq and BS-seq demonstrated the similar sequencing

5    base bias (b). The sequencing quality and data yield are comparable between EM-seq and

6    BS-seq. The final clean data yield was similar at 98%. EM-seq mapping rate was slightly

7    higher. Duplication rate in EM-seq is 50% fewer than in the BS-seq. The global genome

8    coverages were similar between these two methods at 98%. The unmodified cytidine

9    conversion rates were 99.04% in BS-seq and 98.31% in EM-seq. We used the negative control

10    unmodified lambda DNA to calculate the false positive rate (FPR) and the positive control

11    synthetic DNA with specific modified site to calculate the false negative rate (FNR). FPR=C in

12    negative lambda DNA/(C+U in negative lambda DNA);FNR=U in positive PTXB1/(C+U in

13    positive PTXB1) (c). We further look the coverage distribution on Chr4. In the whole genome

14    view, the EM-seq coverage depth are generally higher than the BS-seq in all typical

15    methylation pattern (d). The EM-seq could cover C/CG/CHH/CHG around 25x medium

16    coverage comparing with BS-seq 12x medium coverage(e). To further analyze the relationship

17    between the coverage and GC content. We found both EM-seq and BS-seq tend to cover the

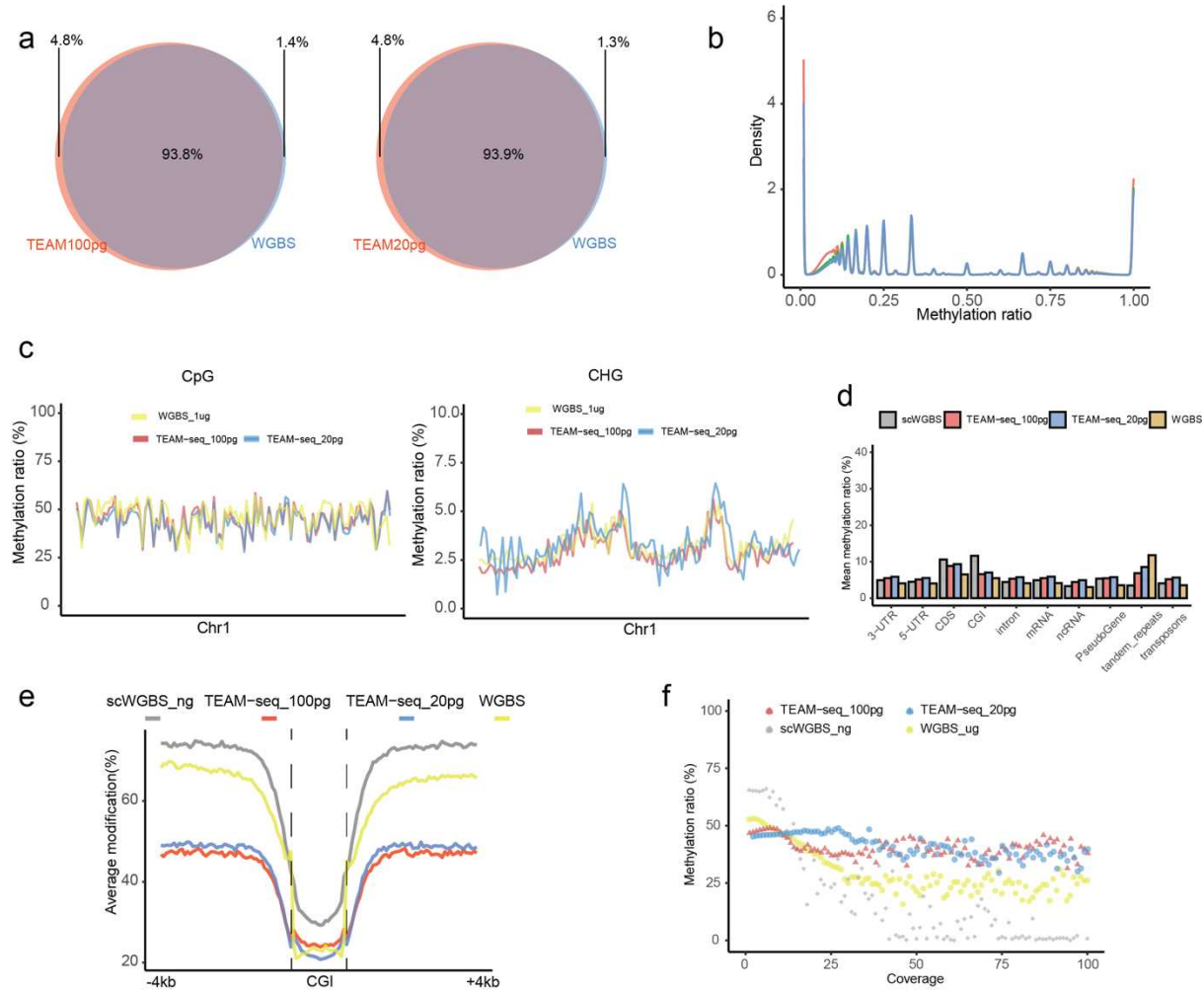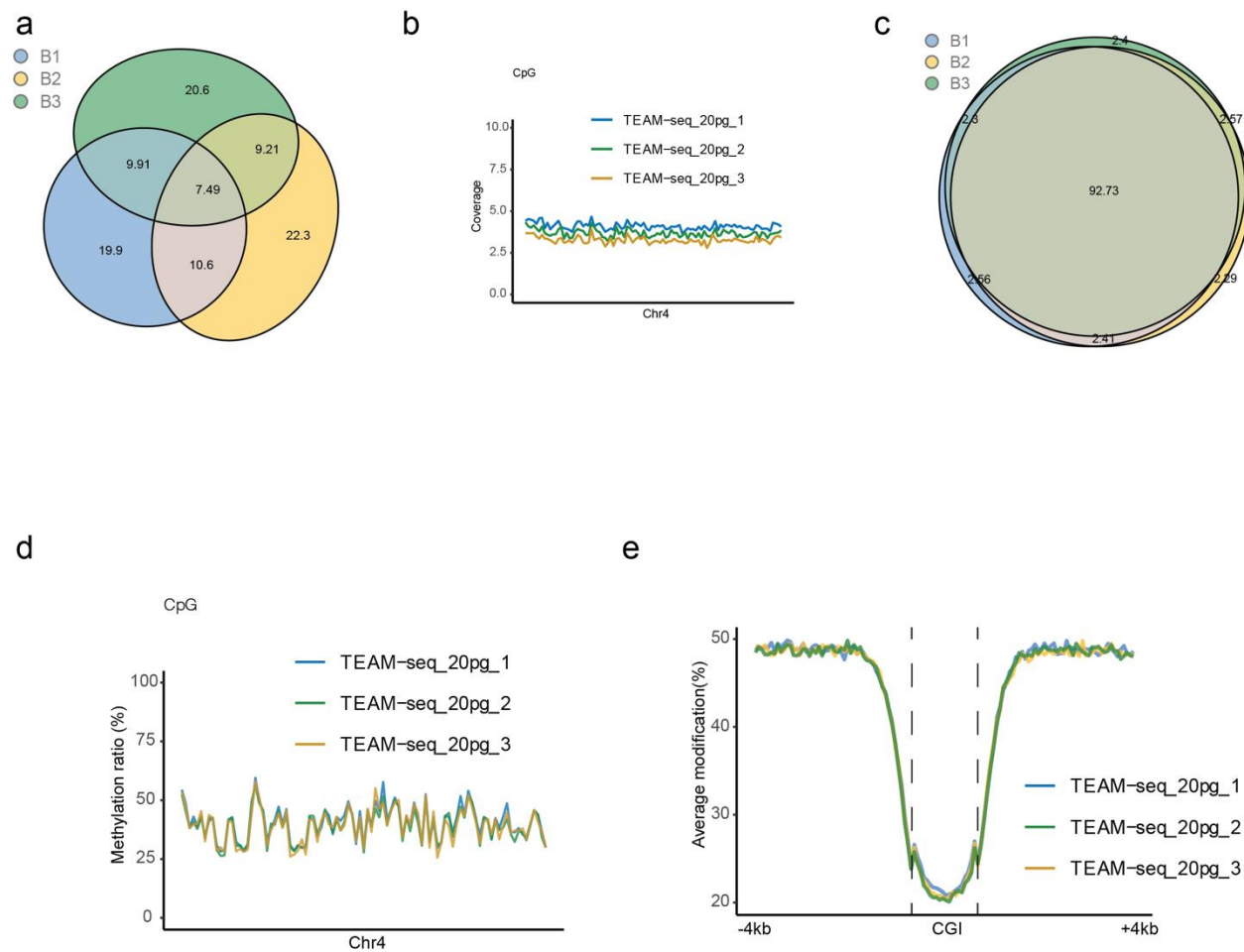18    areas with medium GC content (50%~75%) (f).

19

20

21

a



b



c



d



e



1

1  Sfig2.The sequencing quality and the even coverage of the ultra-low input TEAM-seq.(a) The

2  lines with four colors showed the ATGC proportion in each sequencing base. The TEAM-

3  seq_100pg and TEAM-seq_20pg represented the sequencing result with different initial DNA

4  input 100-pg/20-pg. The 20-pg TEAM-seq showed the more severe fluctuation than the 100-pg

5  TEAM-seq. (b) Left panel showed the accumulation curve of the read covered times per base

6  in 100-pg TEAM-seq and 100-ng scWGBS. The scWGBS used 100-ng as initial DNA input and

7  simplified the experiments to one tube reaction to minimize the loss of operation, and used

8  single stranded library construction [7]. In the TEAM-seq_100pg, 70% sites were covered fewer

9  than 7 times and outperformed the scWGBS_100ng. The scWGBS get to the saturation point

10 (30% genome sites) at 4. We used 20pg genomic DNAs, as the ultra-low input, close to the

11 single cell level. These experiments were repeated for three times. Then we merged these

12 sequencing results of the three 20-pg experimental repeats. The accumulation curve got

13 greatly improvement after merging two 20-pg experimental repeats. (c) The occurrences of

14 covered C/CG/CHG (100 bins) on the chromosome 4 among scWGBS, TEAM-seq, WGBS.

15 The WGBS used 1ug genomic DNA as initial input. All four data were down sampled to the

16 same level. The scWGBS_100ng showed the uneven genomic coverage and some areas

17 were not covered. The WGBS_1ug had the even coverage on the genome with the serious

18 coverage bias on the specific region. Comparing with these two methods, TEAM-seq

19 demonstrated the very even coverage on the Chr4, with ultra-low input DNA. (d) The coverage

20 on the region Chr4:6500000-7000000. TEAM-seq demonstrated the even coverage on this

21 region. (e) The average coverage around the CGI region among scWGBS_100ng,

22 WGBS_1ug, TEAM-seq_100pg, TEAM-seq_20pg. The data were down sampled to the same

1    sequencing depth. The TEAM-seq_20pg demonstrated the highest genome coverage among
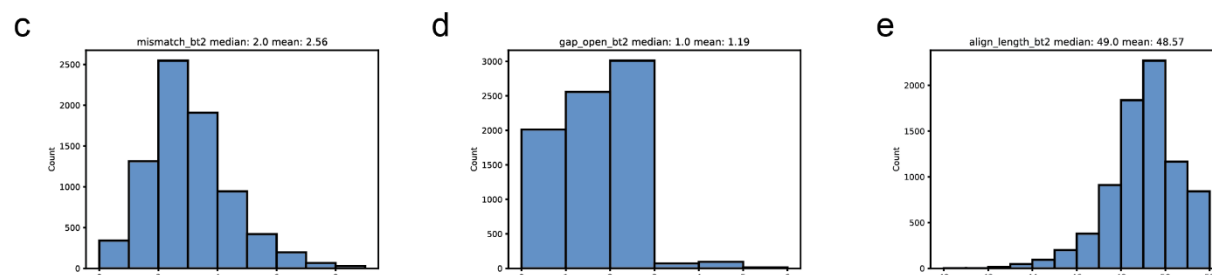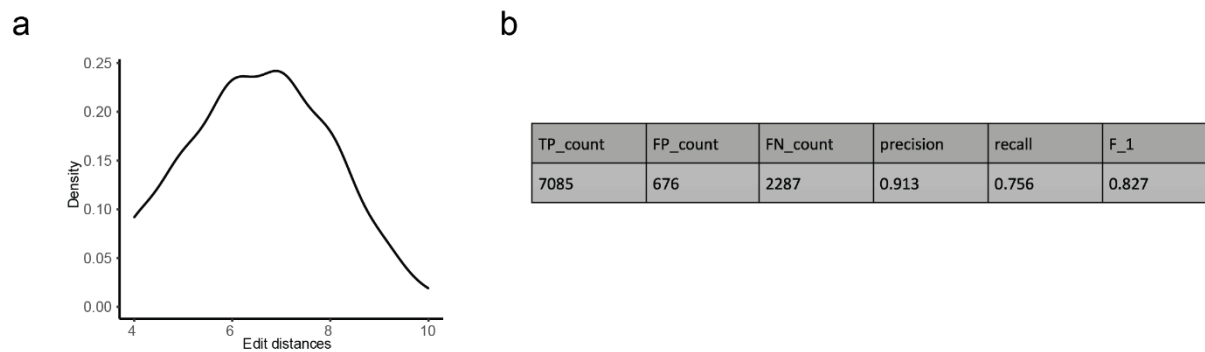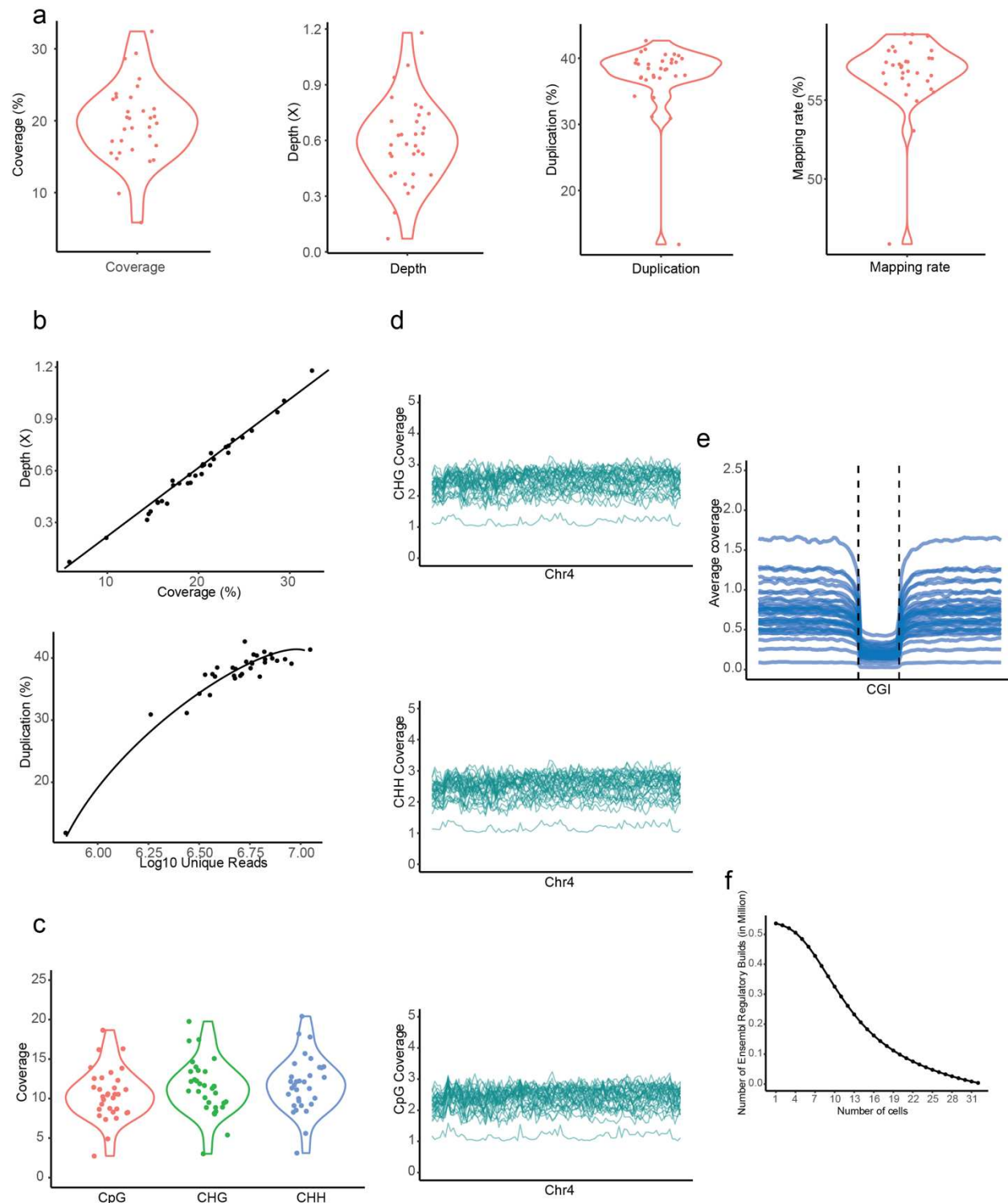
2    these three methods.

3

1    Sfig3. The TEAM-seq and WGBS showed the high consistency in the methylation detection.

2    We selected the high coverage sites (covered by more than 3 reads). The sites with >50%

3    methylation ratio counted as the methylated sites and the sites with <50% methylation ratio

4    counted as the unmethylated sites. The methylation status consistency

5    (methylated/unmethylated) was over 93% between TEAM-seq (20-pg/100-pg) and WGBS

6    (1ug). The calculated methylation ratio correlation between 100-pg TEAM-seq and WGBS was

7    around 0.82 (a). Due to the 20pg genome DNA only contain 1-5 copies of the genomic

8    segments, the methylation ratio was at both of the extremes (0,1) (b). The over methylation

9    level distributions on chr1 (CpG, CHG) were similar between TEAM-seq(20-pg/100-pg) and

10   WGBS(1ug) (c). The TEAM-seq average methylation ratios on the gene element were slightly

11   higher than the WGBS, expect the tandem_repeats (d). We used the CGI as example to study

12   the detailed methylation distribution on the CGI. These four methods all showed the depressed

13   methylation on the CGI and high methylation on the CGI shoulders. The TEAM-seq display

14   30% lower methylation on the CGI shoulders (e). The scWGBS demonstrated the most serious

15   methylation bias resulted by the coverage (f). In the contrast, the TEAM-seq both 100pg and

16   20pg have the very slightly coverage related methylation bias. The results were similar to the

17   previous report TAPS [8].

Sfig4. The high reproducibility in ultra-low input TEAM-seq (20-pg). The 20-pg TEAM-seq was repeated three times independently. The common genome sites, shared by these three independent samples, were 7.49%, due to the sampling bias with minimal amount DNAs (a). The three replicates displayed the similar genome coverage on chr4 (b). The methylation statuses were very consistent between these three replicates (consistency definition as Fig.3). Due to the high methylation consistency, the CpG methylation ratio distribution of these three replicates showed the same pattern on chr4(d). Especially on the CGIs(e), the methylation ratio distributions were nearly the same among these three replicates.
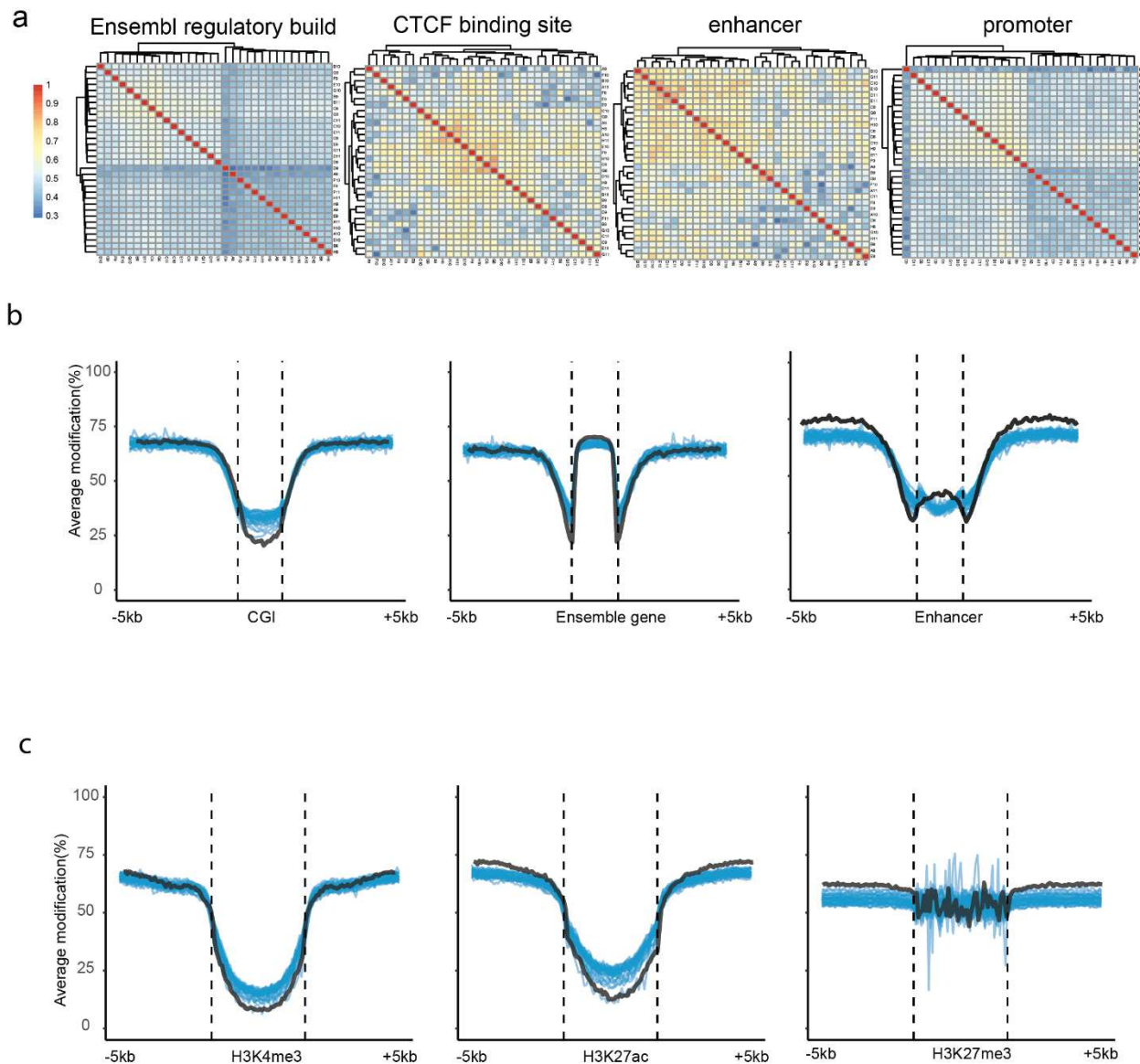
1

2   Sfig5. The barcode designed distance and assessment of the debarcoding methods. We

3   designed the 10 nucleotide barcodes, which own the 4~10 letter differences between each

4   other. The medium edit distance is round 7 (Hamming distance)(a). Then we simulated the

5   nanopore synthetic data, to create enough mismatches, insertions, deletions in barcode and

6   adaptor sequences. We used our house pipeline (method) to debarcode the synthetic data and

7   calculated the TP(True positive), FP(False positive), FN(False negative), precision, recall and

8   F1 (b). To compare the simulated sequence with its original templates, the simulation

9   generated the mean number of mismatches 2.56(c) and the mean number of gaps 1.19(d) on

10  the barcode and adaptor sequence (e, the mean of barcode plus adaptor length 48.57).

11

Sfig6. The coverage of the TEAM-seq on the indexed single cells (HEK293T cells). The medium genome coverage is around ~20% on average ~0.6X sequencing depth. The average duplication rate and mapping rate were ~38% and ~60%(a) (The coverage/depth calculations

1    were described on method). The sequencing depth has the significant effect on the genome

2    coverage with the linear increasing trend, suggesting that the sequencing depth is not close to

3    saturation. The more sequencing data could further improve the genome coverage(b). The

4    duplication rate is hitting the plateau ~38% with 0.6x depth ($10^7$ reads) (b). The mean

5    CpG/CHG/CHH coverage is around ~12%/13%/13%(c). On the whole genome scale, each cell

6    showed the uniform coverage on chr4, which is like 100-pg/20-pg TEAM-seq(d). Moreover,

7    each cell showed the similar coverage pattern on CGI with various sequencing depth(e). Most

8    of the single cell algorithm rely on the common areas covered by most cells, which are

9    important in single cell analysis. In scTEAM-seq (single cell TEAM-seq), ~0.2million Ensembl

10   builds could be covered half cellular population(f).

Sfig7. The accuracy and consistency of the TEAM-seq on the indexed single cells (HEK293T).

The single cell methylation is summarized as the matrix, including the cells (row) and the mean

methylation of ensemble builds. The Ensembl builds carried the location information of CpG

islands, enhancers, promoters, CTCF binding motifs. We then calculated the Pearson

correlations between each cell among these builds. Half of the single cells demonstrated high

correlations in the CTCF binding motifs (~0.7) and enhancers (~0.78) (a). Then we plotted the

methylation distribution pattern of each cell (each blue line represented one cell) on gene

1   elements(b), including CGIs/gene bodies/enhancers, and histone modification motifs(c),

2   including H3K4me3(active)/H3K27ac(active)/H3K27me3(repressive). The black line indicated

3   the methylation distribution of scWGBS (100ng). The methylation distributions were similar

4   between scWGBS and indexed single cell TEAM-seq.

5

6   1.   Smallwood, S.A. et al. Single-cell genome-wide bisulfite sequencing for assessing epigenetic
7     heterogeneity. *Nature Methods* **11**, 817-820 (2014).
8   2.   Miura, F. & Ito, T. Highly sensitive targeted methylome sequencing by post-bisulfite adaptor
9     tagging. *DNA Research* **22**, 13-18 (2015).
10   3.   Kobayashi, H. et al. Repetitive DNA methylome analysis by small-scale and single-cell shotgun
11     bisulfite sequencing. *Genes to Cells* **21**, 1209-1222 (2016).
12   4.   Farlik, M. et al. Single-Cell DNA Methylome Sequencing and Bioinformatic Inference of
13     Epigenomic Cell-State Dynamics. *Cell Reports* **10**, 1386-1397 (2015).
14   5.   Luo, C. et al. Single-cell methylomes identify neuronal subtypes and regulatory elements in
15     mammalian cortex. *Science* **357**, 600-604 (2017).
16   6.   Mulqueen, R.M. et al. Highly scalable generation of DNA methylation profiles in single cells. *Nat*
17     *Biotechnol* **36**, 428-431 (2018).
18   7.   Farlik, M. et al. Single-cell DNA methylome sequencing and bioinformatic inference of
19     epigenomic cell-state dynamics. *Cell Rep* **10**, 1386-1397 (2015).
20   8.   Liu, Y. et al. Bisulfite-free direct detection of 5-methylcytosine and 5-hydroxymethylcytosine at
21     base resolution. *Nature Biotechnology* **37**, 424-429 (2019).
22