

1 **CDCA-seq resolve the different chromatin structure on integral circular DNAs**

2 Weitian Chen<sup>1,2,\*</sup>, Zhe Weng<sup>2,\*</sup>, Zhe Xie<sup>2,3</sup>, Yeming Xie<sup>2</sup>, Chen Zhang<sup>2</sup>, Zhichao Chen<sup>1,2</sup>,  
3 Fengying Ruan<sup>2</sup>, Juan Wang<sup>2</sup>, Yuxin Sun<sup>2</sup>, Yitong Fang<sup>2</sup>, Mei Guo<sup>2</sup>, Yiqin Tong<sup>2</sup>, Yanning Li  
4 <sup>2</sup>, Chong Tang<sup>2,4</sup>

5

6 <sup>1</sup>: BGI Education Center , University of Chinese Academy of Sciences , Shenzhen 518083 ,  
7 China

8 <sup>2</sup>: BGI Genomics, BGI-Shenzhen, Shenzhen 518083, China

9 <sup>3</sup>: Department of Biology, Cell Biology and Physiology, University of Copenhagen 13, 2100

10 Copenhagen , Denmark

11 <sup>4</sup>: Nantong University, Nantong, China, 226000

12 <sup>5</sup>: Nephrosis Precision Medicine Innovation Center, University of Beihua School of Medicine,  
13 Jilin 132011, China

14 \* These authors contributed equally to this work.

15

16 Keywords: ecDNAs, chromatin accessibility, methylation, m<sup>6</sup>A, methyltransferase

17

18 *Running title: comprehensively chromatin accessibility on ecDNAs*

19

20 Correspondence:

21 Chong Tang

22 Director of technology, BGI Shenzhen, China

23 Phone: 8618025420976

24 Email: [tangchong@bgi.com](mailto:tangchong@bgi.com)

25

26

27

28

29

30

1

2 **Abstract**

3 Although ecDNAs have been a subject of sustained research activity for some years, the  
4 underlying mechanism driving the ecDNAs tumorigenesis has begun to unfold recently.

5 Overall, from the results presented in conventional research, the high throughput short reads  
6 sequencing largely ignores the epigenetic status on most ecDNA regions except the

7 junctional areas. We developed a method named CDCA-seq by using methylase to label the  
8 open chromatin without fragmentation, and exonuclease to enrich the ecDNA sequencing

9 depth, followed by the long-read nanopore sequencing. Using this technology, the

10 significantly different patterns of nucleosome/regulator binding were observed in ecDNAs at

11 single-molecule resolution. These results further the understanding of the different

12 regulatory mechanism on ecDNAs.

13

14

15

16

17

18

19

20

21

## 1 **Introduction**

2 Cancer is a malignant disease that is difficult to cure for various reasons, such as oncogene  
3 amplification, tumor evolution, and genetic heterogeneity [1-3]. These have been important  
4 topics of study in the literature for many years. Recently, it has been demonstrated that  
5 extrachromosomal DNA is closely related to carcinogenesis, such as promoting  
6 oncogene amplification[4], driving tumor evolution, and genetic heterogeneity[5-7].  
7 Among extrachromosomal DNAs, extrachromosomal circular DNA (eccDNA), or  
8 microDNA, is the DNA that is arranged next to chromatin in a circular structure,  
9 usually very small in length (< 1 KB). In contrast, circular extrachromosomal DNA  
10 (ecDNA) may be cancer-specific DNA with an average size of 1.3 MB[8]. Although  
11 ecDNA was discovered very early, understanding it has been slow due to early  
12 techniques [9].

13 Many of these techniques are already well-known to study the ecDNA in genetic  
14 and epigenetic fields. Some attempts have been made to solve the ecDNA  
15 identification in sequencing data by improved algorithms. Most algorithms relied on  
16 the detection of ecDNA junction sequence, such as Circle-Map [10], AA  
17 (ampliconarchict) [11], CIRC\_finder[12] , and others [13]. With the aid of these  
18 methods, the ecDNA was identified in numerous cancer tissues [12, 14, 15] [16] [5],  
19 aging cells [17], plasma [18, 19], and healthy somatic tissues [20]. Due to the  
20 rareness of ecDNAs in sequencing data, these approaches require an optimal  
21 solution to enrich the ecDNAs, which obtained circular DNA by digesting linear DNA

1 with nucleases, followed by rolling circle amplification [21]. To further improve the  
2 accuracy of finding ecDNAs, the long-read sequencing technology has also been  
3 used to verify the ecDNA junction structure [22] [23]. Besides the ecDNAs  
4 identification, future advancements are expected to result in the functional  
5 epigenetic study of ecDNAs. Given the rising prevalence of the ecDNAs and  
6 oncogene expression, an essential need is to know the chromatin status and  
7 transcription status of the ecDNAs. The most recent and advanced theory proposed  
8 by Wu et al. offers new insight for the highly accessible chromatin and high  
9 oncogene expression on ecDNAs by ATAC-seq, Chip-seq, and ATACsee [4].  
10 Beyond that, limited studies are available about ecDNA epigenome study because it  
11 is still tricky to analyze its junction structure and epigenome information  
12 simultaneously.

13 Since eccDNAs possess a unique junction sequence pattern, its epigenetic  
14 information can be revealed by locating the neighboring junction regions without  
15 considering distal region, which was thought to be indistinguishable from linear  
16 genome sequences. However, the need for distal coverage of ecDNA epigenome  
17 arose due to the issues with current methods of short reads sequencing and short  
18 fragmentation in ATAC-seq [24] and MNase-seq [25]. Our research is also  
19 motivated by existing methods for assessing chromatin states using long-read  
20 sequencing, such as nanoNOME-seq [26], SMAC-seq [27], fiber-seq[28]. We used  
21 EcoGII, the N6-methyladenosine (m<sup>6</sup>A) methyltransferase, to soft label chromatin

1 accessible regions without fragmentation (Circular DNA Chromatin Accessibility  
2 Sequencing, CDCA-seq). We enriched the ecDNAs by using nuclease to digest the  
3 linear genome. The nanopore sequencing accurately detected the m<sup>6</sup>A probed  
4 chromatin accessible regions and junctional structure properties simultaneously on  
5 ecDNAs in a long-range. Using this method, we found the high diversity of ecDNAs  
6 chromatin accessibility and coordination with distal regulators in single-molecule  
7 resolution, which has not been reported before.

8

## 9 **Results**

### 10 **CDCA-seq comprehensively maps chromatin accessibility and nucleosome** 11 **positioning in ecDNAs at multikilobase scale**

12 The ecDNAs play a vital role in tumorigenesis due to their high chromatin  
13 accessibility and amplified oncogene expression [4]. Conventional approaches to  
14 study chromatin accessibility are based on the concept that the chromatins  
15 protected the bound sequence without being attacked by a transposon (Figure 1A)  
16 or MNase [25]. In the ATAC-seq, the transposon preferably tagmented the openly  
17 accessible genome region, followed by the next-generation sequencing (NGS)  
18 (Figure 1A). However, these are not employed in most integral ecDNA chromatin  
19 studies due to the homologous ecDNA/genome sequences, making the problematic  
20 distinction between ecDNA and linear genome. In general, prior work in ecDNA

1 chromatin studies, based on NGS short reads, only observed the chromatin status  
2 in the junction region (200bp around junction) and took minimal consideration of the  
3 other distal areas on ecDNAs owing to limitation of techniques (>200bp to junction  
4 regions) (Figure 1A). We built a generalized framework for solving problems in  
5 the integral ecDNA chromatin studies on the concept of the SMAC-seq[27] and  
6 fiber-seq[28]. We applied the soft labeling with EcoGII, the m<sup>6</sup>A methyltransferase,  
7 preferably methylating the adenosine on the openly accessible DNA region without  
8 fragmentation by a transposon (Figure 1A). To improve the ecDNA capturing  
9 efficiency, the exonuclease was introduced to remove the linear genome DNA[29].  
10 The integral ecDNAs could be sequenced by nanopore, and the probed m<sup>6</sup>A could  
11 be detected [27]. In the following data analysis, we first identified the ecDNAs by  
12 head-to-tail junction locations with dynamically mapping the segments of  
13 sequences to the genome (Figure 1B). Based on the head-to-tail junction locations,  
14 we then reassembled the partial ecDNA sequences as the new reference and  
15 identify the m<sup>6</sup>A signal based on the reassembled ecDNA sequence to prevent  
16 signal bias occurring in the junction region (Figure 1B).

17 Fundamental read length statistical analysis showed that the mean read length is  
18 around 10kb up to 100kb, which is 50x broader than the junctional region observed  
19 in conventional ATAC-seq[4] (Supplemental Figure 1). The long-read feature also  
20 makes the nanopore method optimal for applications such as SV, CNV, ecDNA  
21 identification with better sensitivity and specificity [30]. As expected, 80% ecDNAs

1 detected in our CDCA-seq could be validated through PCR (Supplemental Figure 2).  
2 ecDNAs and residual linear DNAs account for 0.9% and 99.1% in the total  
3 sequencing reads (Supplemental Figure 3) with exonuclease treatment. The m<sup>6</sup>A  
4 possibility distribution in Megalodon, showed two distinct peaks for the treated sample. The  
5 distribution of the narrow peak with lower m<sup>6</sup>A probability (mean=0.49) is similar to the  
6 background noise distribution (Supplemental Figure 4). Therefore, we set m<sup>6</sup>A  
7 methylation probability over 0.53 as the cutoff for true m<sup>6</sup>A signal. (Supplemental  
8 Figure 4). The real positive cutoff value was set as 0.53, and the m<sup>6</sup>A calling  
9 specificity and sensitivity is 0.99 and 0.92 (Supplemental Figure 4). The residual  
10 linear DNA was used as an internal control to validate against the published  
11 ATAC-seq data[31]. The CDCA-seq achieved consistency and coherence with  
12 ATAC-seq data in various resolutions (Figure 1C, Supplemental Figure 5). Good  
13 concordance was also found when comparing our results against the published  
14 methods [27, 28]. The m<sup>6</sup>A labeling deviation has been strongly reduced to 0.015  
15 when accounting for the m<sup>6</sup>A mean ratio (Supplemental Figure 6). The impact of  
16 exonuclease treatment and reproducibility has been also investigated here.  
17 (Supplemental Figure 7). These characteristics of CDCA-seq are critical for  
18 effectively measuring the chromatin accessibility on linear and circular DNAs in the  
19 multikilobase range.  
20 Another remarkable feature of CDCA-seq is the single molecular resolution  
21 chromatin status on ecDNAs. At the single molecular level, the single base m<sup>6</sup>A

1 possibility varied from 0.6 to 1 other than the genome level, whereas the genomic  
2 level, the m<sup>6</sup>A possibility could be improved by multiple covered reads  
3 (Supplemental Figure 8). In practice, the resolution of chromatin accessibility is  
4 around 200bp. We adopted a Bayesian procedure to aggregate methylation  
5 probabilities and derived the accurate single-molecule accessibility calls over  
6 arbitrary size windows (Supplemental Figure 8). In summary, the CDCA-seq offers  
7 attractive features in terms of integral ecDNA chromatin status in the multikilobase  
8 range at the single molecular resolution.

9

#### 10 **The diversified chromatin accessibility on ecDNAs**

11 Evidence from other studies by ATAC-seq and Chip-seq suggests that the active  
12 chromatin status and highly accessible chromatin on ecDNAs may be associated  
13 with the high transcription level of oncogenes [4]. To distinguish the ecDNAs  
14 molecules with linear DNA molecules in ATAC-seq and Chip-seq, we first need to  
15 screen out the short reads (~200bp) spanning the nonhomologous end-joining  
16 sequence of ecDNAs. One problem with these approaches is the potential bias to  
17 neglect the distal regions due to focusing on the reads of (~200bp) neighboring  
18 junctional sequence of ecDNAs. The long-read technology CDCA-seq may offer  
19 advantages for precisely detecting ecDNAs [30] [32] [22] and observing the distal  
20 chromatin status on integral ecDNAs. We observed an extensive catalog of 12997  
21 different ecDNAs formed from chromosomal breakpoints between 0.05kb and up to



1 100kb (Supplemental Table 1). Gene ontology (GO) analysis of these ecDNA  
2 carried genes identified some significantly enriched GO terms, including  
3 GTPase-related activity, channel activity, nucleoside-triphosphatase activity, which  
4 play essential roles in cancer progression [33] (Supplemental Figure 9). By  
5 RNA-seq data analysis[34], the highly expressed ecDNA carried genes (25% rank),  
6 medium expressed genes (25~75% rank), and lowly expressed genes (75%~100%  
7 rank) are counted as 340, 464, and 589 respectively, indicating that not all the  
8 ecDNA carried genes are highly expressed.

9 By comparing the average chromatin accessibility between the ecDNAs and  
10 homologous linear DNAs, we found that the overall chromatin status of ecDNAs is  
11 2x more accessible than that of linear DNAs (Figure 2A). These findings reinforce  
12 the general belief that the ecDNA amplification resulted in higher oncogene  
13 transcription [4], coupled with the enhanced chromatin accessibility in the junctional  
14 region. The data were subjected to the detailed mapping of the ecDNA chromatin  
15 status. We found that the chromatins in the ecDNA junctional areas are significantly  
16 more accessible than in other linear homologous regions (Figure 2B). This is an  
17 interesting find, as it suggests that the conclusion drawn by only observing the  
18 junctional areas of conventional ATAC-seq may be bias for the whole ecDNA  
19 chromatin. We calculate the average m<sup>6</sup>A methylation fraction covering from the  
20 TSS (gene transcription start site) to the TES (gene transcription end site) on each  
21 gene-spanning read. A pairwise scatter plot of the average accessibility between

1 ecDNAs carried genes and linear genome had genes showed that 63% of the gene  
2 regions are more accessible on ecDNAs than linear DNAs (Figure 2C). Comparing  
3 the ecDNA and linear DNA chromatin profiles around the TSS/TES (+/- 500)  
4 revealed a significant difference in nucleosome depletion/occupancy patterns  
5 (Figure 2DE). The order nucleosome organization may impact access to ecDNA  
6 (Figure 2DE). Considering that the 63% gene regions are more accessible on the  
7 ecDNAs than on the linear DNA, we further plotted the chromatin structure around  
8 TSS/TES (+/-500bp) of these genes (Supplemental Figure 10). The formation of  
9 nucleosome depletion regions (NDRs) on linear DNAs is restricted to 200bp ahead  
10 TSSs. In contrast, the NDRs on ecDNAs are distributed in order (Supplemental  
11 Figure 10). The other 37% gene regions are more accessible on linear DNAs than  
12 ecDNAs. The TSSs/TESSs (+/- 500bp) are also significantly more accessible on  
13 linear DNAs than ecDNAs with different NDR patterns (Supplemental Figure 11).  
14 The formation of large NDRs is restricted to TSSs on linear DNAs, which is not  
15 observed on ecDNAs.

16 Another illustration of the complex interplay between chromatin states on ecDNAs  
17 and linear DNAs relates to the transcriptional activity. The linear DNA chromatin  
18 state on the active genes (top 25% rank) largely devoid of nucleosomes on the  
19 TSSs due to extremely high transcription activity (Supplemental Figure 12). In  
20 contrast, the ecDNA chromatin structure of active genes adopts the distinct  
21 conformation, implying the different regulatory mechanisms on ecDNAs

1 (Supplemental Figure 12). For the transcriptionally inactive genes, the stationary  
2 nucleosome states are shown on the linear DNAs (Supplemental Figure 13). In  
3 contrast, the ecDNAs still have the active nucleosome organization on 300bp ahead  
4 TSSs, suggesting that chromatin accessibility is necessary but not sufficient for  
5 enhancer or promoter activity on ecDNAs (Supplemental Figure 13). In conclusion,  
6 the ecDNAs and linear DNAs have a significantly distinct nucleosome  
7 depletion/occupancy pattern in various conditions, suggesting the different gene  
8 regulatory mechanisms on between ecDNAs and linear DNAs.

### 9 **The chromatin status on ecDNA and linear genome on the single-molecular** 10 **resolution**

11 The conventional ATAC-seq is based on statically calling the peak of the enriched  
12 read in a specific region[24]. Recent single-molecule and single-cell measurements  
13 of accessibility suggest that ATAC-seq on cell populations represent an ensemble  
14 average of distinct molecular states [35]. An essential attribute of the CDCA-seq  
15 measures the ecDNA chromatin accessibility in single molecular resolution by  
16 taking the advantages of the small variance (Supplemental Figure 6) and increased  
17 cumulative possibility in segments (Supplemental Figure 8). Measuring chromatin  
18 accessibility of the single linear DNA has also been done in the SMAC-seq [27] and  
19 fiber-seq [28].

20 We then asked whether the CDCA-seq could reveal multiple chromatin accessibility  
21 states on ecDNAs. The chromatin structure on linear DNAs

1 (chr10:42383201~42389251) adopts two distinct conformations; an inactive  
2 nucleosomal state and a state largely devoid of nucleosomes due to extremely high  
3 transcription activity [36] (Figure 3A). It was conventionally believed that the active  
4 nucleosome status occupied the majority in cancer cells. As expected, the 70%  
5 proportion of ecDNAs molecule comes from the very active chromatin state (Figure  
6 3A). We observed the highly heterogeneous nucleosome depletion/occupancy  
7 pattern on ecDNAs, and most chromatin molecules are not very active in the  
8 positive strand, suggesting the different transcription regulation on ecDNAs (Figure  
9 3B upper panel). The different phenomena were also observed in other regions  
10 (Supplemental Figure 14). To avoid the bias conclusion drawn by the methylase  
11 heterogeneous activity, the other upstream and downstream region is chosen as  
12 quality control.

13 To further explore the limit of CDCA-seq's resolution limits, we studied methylation  
14 patterns in more detail. We next quantified strand-specific DNA accessibility and  
15 observed a strand-asymmetric DNA accessibility pattern in the linear genome  
16 (Supplemental Figure 14). The strand-asymmetric DNA accessibility pattern is also  
17 observed in ecDNAs, and both strands display high heterogeneity (Figure 3B,  
18 supplemental Figure 14). This strand-specific heterogeneity in methylation potential  
19 within the nucleosome may inform how transcription factors might interact with  
20 nucleosome-associated DNA in vivo.

1 Wu et al. show ecDNA enables ultra-long-range chromatin contact, permitting distant  
2 interactions with regulatory elements [4]. We next examined co-accessibility patterns in the  
3 ecDNAs and linear genome by assessing nucleosome positioning correlations. The  
4 nucleosomes have a higher correlation on ecDNAs than linear DNAs (Figure 4AB,  
5 Supplemental Figure 15). Moreover, the ecDNAs and linear DNAs adopt significantly  
6 different chromatin co-accessibility patterns (Figure 4AB, Supplemental Figure 15). Average  
7 co-accessibility profiles on linear DNAs reveal the detectable correlation between  
8 nucleosome positions up to two to three nucleosomes away. The ecDNAs correlation may be  
9 further and up to 20 nucleosomes away (Figure 4AB, Supplemental Figure 15). These results  
10 agree with the HiC result [4], that the ecDNAs have the distant chromatin interaction. It was  
11 interesting to note that the ecDNA demonstrates some ultra-distant anticorrelated states.  
12 Overall, the ecDNAs have high heterogeneity in each molecule and remote chromatin  
13 interaction, suggesting the different regulation mechanisms from linear DNAs.

## 14 **Discussion**

15 An understanding of ecDNAs may prove to be essential for tumorigenesis [5-7]. Most of the  
16 community's efforts have gone into solving ecDNA identification in various cancer tissues [12,  
17 14, 15] [16] [5]. There has been an increasing research effort, and specialization in ecDNA  
18 open chromatin status to resolve the oncogene amplification on ecDNAs[4]. However, most  
19 studies focused on the short sequencing reads with junctional sequences detected to avoid  
20 the false-positive identification of ecDNAs and precisely determine the ecDNA epigenetic

1 status. A large subgroup (38%) and a large part (60%) of the ecDNAs covered regions of  
2 DNAs that are not unique in the reference genome, which is difficult to be discriminated [37].  
3 In this study, we use nanopore sequencing to evaluate integral ecDNA chromatin  
4 accessibility on long strands of ecDNAs by applying m<sup>6</sup>A methyltransferase to label open  
5 chromatin without fragmentation. Consistent with the findings of the other [4], 63% of the  
6 ecDNA carried genes have more accessible chromatin structure than the linear DNA.  
7 However, other 37% ecDNA chromatin on gene regions are less accessible than  
8 corresponding linear DNAs. Notably, the nucleosome depletion/occupancy patterns are  
9 significantly different between ecDNAs and linear DNAs. This single-molecule resolution  
10 method allows footprinting of protein and nucleosome binding and determination of  
11 epigenetic signature on chromatin accessibility. It is hoped that this study will lead to the new  
12 insight of comprehensively understanding ecDNA epigenome regulation.

13 Move to result section, our study would treat sample DNAs with an exonuclease  
14 that would remove most of the linear DNA molecules in the sample and increase the  
15 sequencing depth for the ecDNAs (0.9%). Few identified linear DNAs maybe come  
16 from ecDNA homologous regions without junctions, and the likelihood is around  
17 0.9%, which is neglectable. Comparing with non-digest direct sequencing, we only  
18 get 0.1% ecDNA related reads (Supplemental Figure 16). The circular eccDNA  
19 enrichment fold is 10x. The exonuclease treatment not only improve the ecDNA  
20 sequencing coverage, but also the ecDNA detection specificity (Supplemental  
21 Figure 2). However, the DNA purification process would damage large-sized DNAs,

1 especially for the ecDNAs larger than 1MB [38]. The damaged ecDNAs could be  
2 digested during exonuclease digestion and missed in the sequencing. A method  
3 that gently purifies the large DNAs is preferable in further mega ecDNA studies.

4 A newer software Megalodon of m<sup>6</sup>A signal calling is chosen and compared with Tombo,  
5 which is presented in other literature [27] [28]. In the ecDNAs m<sup>6</sup>A calling, the Tombo would  
6 ignore half of the sequence or lose most ecDNAs for unknown reasons (Supplemental Figure  
7 17). The Tombo lost 83% sensitivity on the ecDNA m<sup>6</sup>A signal calling. Although the  
8 Megalodon improves the sensitivity on ecDNAs m<sup>6</sup>A calling, it does not address the issue of  
9 the false-positive m<sup>6</sup>A signal. Most adenosine could be recognized as m<sup>6</sup>A with the possibility  
10 of 0.4~1 in Megalodon, other than the Tombo, which has higher specificity. The only known  
11 way to solve this problem is utilizing data training with negative control samples  
12 (Supplementary Figure 4). We used 0.53 as m<sup>6</sup>A possibility cutoff, successfully discriminating  
13 the m<sup>6</sup>A and false-positive with sensitivity 0.92 and specificity 0.99. In general, the  
14 Megalodon has demonstrated better performance in ecDNA research and as well as  
15 improved the specificity with our data training.

16 In the sequencing data, we found that the methylated treated DNA generated more  
17 data than the non-methylated DNAs, which is not consistent with  
18 SMAC-seq/fiber-seq[27] [28]. The highly open chromatin with high methylated sites  
19 may be enriched in our method. In our lab experiments, we found the heavily  
20 modified DNA is more resistant to exonuclease digestion, leading to the enrichment

1 of modified DNAs. The non-treated sample showed the lower overall methylation  
2 level (Supplemental Figure 18). However, the nucleosome occupancy positions  
3 were not significant affected by the exonuclease treatment (Supplemental Figure  
4 19). Moreover, in the strand-specific view, the reverse strand reads are generally  
5 less abundant than the positive strands. This may also be due to the different  
6 methylation statuses on the positive and negative strands, which lead to different  
7 digestion effectiveness. This problem is usually overcome by increasing  
8 sequencing depth and normalization methods. We also suggest sequencing both  
9 treated and non-treated samples for ecDNA sequencing coverage and further  
10 improved quantification accuracy.

11 Only 63% of highly chromatin-accessible gene regions are observed in our  
12 experiment. However, Wu, etc., showed that the ecDNAs are overall highly  
13 chromatin accessible by ATACsee technology [4]. When comparing results from all  
14 other regions against the published data, good agreement was found that 80% of  
15 areas are highly accessible on the ecDNAs (Supplemental Figure 20). Most of the  
16 highly chromatin-accessible areas are distributed in the intron and intergenic  
17 regions (Supplemental Figure 21). The reasons for this remain unclear, but our  
18 results indicate that the ecDNAs have a highly open chromatin structure, especially  
19 in the intergenic and intronic regions.



1 The CDCA-seq provides the useful functionality to study the chromatin status on the integral  
2 ecDNAs, offering a deep insight into the different regulation mechanism of ecDNAs. The  
3 ecDNA enrichment step is to resort to using exonuclease treatment, causing the loss of mega  
4 ecDNAs. It is assumed that future advances will help to address these problems by DNA  
5 damages in the purification and sequencing depth.

6

## 7 **Method**

### 8 Cell culture

9 Human mammary gland carcinoma cell line MCF-7 was obtained from ATCC. MCF-7 were  
10 grown in DMEM(Gibco,11995065) supplemented with 10% FBS (Gibco,10099141),  
11 0.01mg/ml insulin(), and 1% penicillin-streptomycin(Gibco, 15140122). The cell line was  
12 regularly checked for mycoplasma infection (Yeasen, 40612ES25).

### 13 Nuclei isolation and MTase treatment

14 Cells were grown to 70-80% confluency, and were collected by TrypLE (Gibco,12604013).  
15 After 300xg centrifuge for 5 minutes, nuclei were isolated with lysis buffer (100 mM Tris-HCl,  
16 pH 7.4, 10 mM NaCl, 3 mM MgCl<sub>2</sub>, 0.1 mM EDTA, 0.5% CA630) for 5 minutes on ice. Nuclei  
17 were centrifuged at 300xg in wash buffer (100 mM Tris-HCl, pH 7.4, 10 mM NaCl, 3 mM  
18 MgCl<sub>2</sub>, 0.1 mM EDTA) at 4 degree, and washed twice for 5 minutes and counted.

19 1x10<sup>6</sup> intact nuclei were subjected to an m<sup>6</sup>A methylation reaction mixture containing 1x  
20 Cutsmart buffer (NEB), 200U of non-specific adenine methyltransferase M.EcoGII (NEB,

1 M0603S), 300mM sucrose, and 96  $\mu$ M S-adenosylmethionine (NEB, B9003S) in 500ul  
2 volume. The reaction mixture was set up at a 37-degree thermomixer with shaking at  
3 1000rpm for 30 minutes. S-adenosylmethionine was replenished at 640uM every 7.5  
4 minutes at 7.5, 15, and 22.5 minutes into the reaction mixture. The reaction was stopped by  
5 adding an equal volume of stop buffer (20 mM Tris-HCl pH 7.4, 600 mM NaCl, 1% SDS, 10  
6 mM EDTA). No methylation controls were treated in the same conditions without adding  
7 M.EcoGII in the reaction mixture. The samples were then treated with 20ul of Proteinase K  
8 (20mg/ml) at 55 degrees overnight, and the DNA was extracted with phenol: chloroform  
9 extraction and ethanol precipitation.

10 ecDNA isolation, purification, and sequencing

11 ecDNA was isolated by Circle-Seq[21] method, which digested linear DNA with  
12 modifications. Briefly, 10ug of M.EcoGII treated DNA was subjected to a reaction mixture  
13 containing 1x plasmid-safe reaction buffer, 20U plasmid-safe ATP-dependent DNase  
14 (Lucigen, E3101K), 1mM ATP, and nuclease-free water was supplemented to a final volume  
15 of 100ul. The reaction mixture was incubated at 37 degrees for 7 days. Every 24 hours, the  
16 reaction mixture was replenished by adding 20U plasmid-safe ATP-dependent DNase, 1mM  
17 ATP, and 0.4ul 10X plasmid-safe reaction buffer. Digested ecDNA was purified with 1.8X  
18 Ampure XP beads (Beckman Coulter).

19 Purified ecDNA was prepared for nanopore sequencing by ligation kit LSK-SQK108(ONT).

20 The samples were 10kb by Covaris G tubes, end-repaired and dA-tailed using NEBnext

21 Ultra II end-repair module (NEB), followed by clean-up using 1.8X Ampure XP beads.

1 Sequencing adaptors and motor proteins were ligated to end-repaired DNA fragments using  
2 blunt/TA ligase master mix (NEB), followed by clean-up using 0.4x AMPure XP beads. 1ug  
3 adaptor-ligated samples per flow cell were loaded onto PRO-002 flowcells and run on  
4 PromethION sequencers for up to 72h. Data were collected by MinKNOW v.1.14.

#### 5 Base-calling and Linear DNA methylation calling

6 Reads from the ONT data were processed using Megalodon V2.2.9, which used Guppy  
7 base caller to base-calling, and Guppy model config  
8 `res_dna_r941_min_modbases-all-context_v001.cfg` released into the Rerio repository was  
9 used to identify DNA m<sup>6</sup>A methylation. `Megalodon_extras` was used to get per read  
10 `modified_bases` from the Megalodon basecalls and mappings results. To further explore the  
11 accurate threshold of methylation probability, a control sample with almost no m<sup>6</sup>A  
12 methylation was used as background noise, and the Gaussian mixture model was used to fit  
13 the methylation probability distribution generated by Megalodon.

#### 14 ecDNA calling

15 ONT Reads meet the following conditions were defined as ecDNA molecules performed by  
16 the inner mappy/minimap2 aligner [39]. (1) One segment (>1kb) of an ONT read was  
17 mapped to the genome at one site, and another segment (>1kb) was mapped to the genome  
18 at another site. (2) Two segments were mapped to the same chromosome. (3) Two  
19 segments were mapped to the same strand of the genome. (4) Two segments in a pair  
20 showed outward orientation.

## 1 Nanopore ecDNA methylation calling

2 Due to ecDNA special structure, the m<sup>6</sup>A calling cannot be successfully performed by  
3 aligning to the reference genome, especially for junctional regions. The custom python script  
4 was used to assemble ecDNA reference genome sequences according to the table  
5 generated from the previous step. Considering that the read length might be longer than the  
6 ecDNA reference, the ecDNA reference was subsequently preprocessed by adding 10M N  
7 to the ends to increase the mapping efficiency. The downstream step is performed in a  
8 similar way as linear DNA methylation calling.

## 9 Annotation and methylation configuration

10 TES, TTS, CDS, and other gene elements were downloaded from UCSC Table Browser,  
11 And the gene elements were processed into 50bp bin for downstream analysis. Linear DNA  
12 and ecDNA were also processed to the size of 50bp bin and sliding for 5bp. The accessibility  
13 score over multi base-pair windows was calculated as methylation ratio = m<sup>6</sup>A bases in all  
14 covered reads under bin/ adenosine bases in all covered reads under the bin.

## 15 RNA-seq data analysis

16 The RNA-seq data of MCF-7 was downloaded from the Gene Expression Omnibus (GEO)  
17 repository database with the accession number GSE71862. The gene expression was  
18 divided into three categories: high, medium, and low, representing 25%, 25%-75%, and 75%  
19 gene expression rank, respectively.

## 1 Co-accessibility assessment

2 To evaluate co-accessibility patterns along the genome, we applied COA as follows. Each  
3 chromosome in the genome was split into windows of size  $w$ . For each such window  $(i, i +$   
4  $w)$ , we identified another window  $(j, j+w)$  such that the span  $(i, j, w)$  was covered by  $\geq N$  reads.  
5 For each single spanning molecule  $k$ , accessibility scores ( $A$ ) in each bin were then  
6 aggregated and binarized as described above. The local co-accessibility matrix between two  
7 windows was calculated:

$$8 \quad COA_{i,j,w} = \text{mean}_N \left( 1 - \frac{|A_{i,w} - A_{j,w}|}{A_{i,w} + A_{j,w}} \right)$$

## 9 Data availability

10 The data that support the findings of this study have been deposited into CNGB Sequence  
11 Archive (CNSA) [40] of China National GeneBank DataBase (CNGBdb)[41] with accession  
12 number CNP0001299.

## 13 **Acknowledgment**

14 Funding: This research was supported by the Science, Technology, and Innovation  
15 Commission of Shenzhen Municipality (grant number JSGG20170824152728492). The  
16 supporter had no role in designing the study, data collection, analysis, and interpretation, or  
17 in writing the manuscript.

## 18 **Author contributions**

1 C.T. designed and supervised the experiments. Z.W. and X.Z perform the lab experiments;  
2 W.T.C. performs the bioinformatics data analysis. All others joined the data analysis.

### 3 **Competing interest**

4 The authors declare no competing interests.

5

### 6 **References**

- 7 1. Gillies RJ, Verduzco D, Gatenby RA: Evolutionary dynamics of carcinogenesis and why targeted  
8 therapy does not work. *Nat Rev Cancer* 2012, 12:487-493.
- 9 2. Ray Chaudhuri A, Callen E, Ding X, Gogola E, Duarte AA, Lee JE, Wong N, Lafarga V, Calvo JA,  
10 Panzarino NJ, et al: Replication fork stability confers chemoresistance in BRCA-deficient cells. *Nature*  
11 2016, 535:382-387.
- 12 3. Turajlic S, Swanton C: Implications of cancer evolution for drug development. *Nat Rev Drug Discov*  
13 2017, 16:441-442.
- 14 4. Wu S, Turner KM, Nguyen N, Raviram R, Erb M, Santini J, Luebeck J, Rajkumar U, Diao Y, Li B, et al:  
15 Circular ecDNA promotes accessible chromatin and high oncogene expression. *Nature* 2019,  
16 575:699-703.
- 17 5. Turner KM, Deshpande V, Beyter D, Koga T, Rusert J, Lee C, Li B, Arden K, Ren B, Nathanson DA, et  
18 al: Extrachromosomal oncogene amplification drives tumour evolution and genetic heterogeneity.  
19 *Nature* 2017, 543:122-125.
- 20 6. Paulsen T, Kumar P, Koseoglu MM, Dutta A: Discoveries of Extrachromosomal Circles of DNA in  
21 Normal and Tumor Cells. *Trends Genet* 2018, 34:270-278.

- 1 7. Verhaak RGW, Bafna V, Mischel PS: Extrachromosomal oncogene amplification in tumour  
2 pathogenesis and evolution. *Nature Reviews Cancer* 2019, 19:283-288 %@ 1474-1768.
- 3 8. Chiu RWK, Dutta A, Henssen AG, Lo YMD, Mischel P, Regenber B: What is extrachromosomal  
4 circular DNA and what does it do? *Clinical Chemistry* 2020, 66:754-759 %@ 0009-9147.
- 5 9. Bailey C, Shoura MJ, Mischel PS, Swanton C: Extrachromosomal DNA-relieving heredity constraints,  
6 accelerating tumour evolution. *Ann Oncol* 2020, 31:884-893.
- 7 10. Prada-Luengo I, Krogh A, Maretty L, Regenber B: Sensitive detection of circular DNAs at  
8 single-nucleotide resolution using guided realignment of partially aligned reads. *BMC Bioinformatics*  
9 2019, 20:663.
- 10 11. Deshpande V, Luebeck J, Nguyen N-PD, Bakhtiari M, Turner KM, Schwab R, Carter H, Mischel PS,  
11 Bafna V: Exploring the landscape of focal amplifications in cancer using AmpliconArchitect. *Nature*  
12 *communications* 2019, 10:1-14 %@ 2041-1723.
- 13 12. Kumar P, Kiran S, Saha S, Su Z, Paulsen T, Chatrath A, Shibata Y, Shibata E, Dutta A: ATAC-seq  
14 identifies thousands of extrachromosomal circular DNA in cancer and cell lines. *Science Advances*  
15 2020, 6:eaba2489 %@ 2375-2548.
- 16 13. Chen FZ, You LJ, Yang F, Wang LN, Guo XQ, Gao F, Hua C, Tan C, Fang L, Shan RQ: CNGBdb:  
17 China National GeneBank DataBase. *Yi Chuan= Hereditas* 2020, 42:799-809 %@ 0253-9772.
- 18 14. Koche RP, Rodriguez-Fos E, Helmsauer K, Burkert M, MacArthur IC, Maag J, Chamorro R,  
19 Munoz-Perez N, Puiggròs M, Garcia HD: Extrachromosomal circular DNA drives oncogenic genome  
20 remodeling in neuroblastoma. *Nature Genetics* 2020, 52:29-34 %@ 1546-1718.
- 21 15. Paulsen T, Shibata Y, Kumar P, Dillon L, Dutta A: Extrachromosomal circular DNA, microDNA, without

- 1 canonical promoters produce short regulatory RNAs that suppress gene expression. *bioRxiv*  
2 2019:535831.
- 3 16. Verhaak RGW, Bafna V, Mischel PS: Extrachromosomal oncogene amplification in tumour  
4 pathogenesis and evolution. *Nat Rev Cancer* 2019, 19:283-288.
- 5 17. Hull RM, King M, Pizza G, Krueger F, Vergara X, Houseley J: Transcription-induced formation of  
6 extrachromosomal DNA during yeast ageing. *PLoS Biol* 2019, 17:e3000471.
- 7 18. Zhu J, Zhang F, Du M, Zhang P, Fu S, Wang L: Molecular characterization of cell-free eccDNAs in  
8 human plasma. *Sci Rep* 2017, 7:10968.
- 9 19. Kumar P, Dillon LW, Shibata Y, Jazaeri AA, Jones DR, Dutta A: Normal and Cancerous Tissues  
10 Release Extrachromosomal Circular DNA (eccDNA) into the Circulation. *Mol Cancer Res* 2017,  
11 15:1197-1205.
- 12 20. Møller HD, Mohiyuddin M, Prada-Luengo I, Sailani MR, Halling JF, Plomgaard P, Maretty L, Hansen  
13 AJ, Snyder MP, Pilegaard H, et al: Circular DNA elements of chromosomal origin are common in  
14 healthy human somatic tissue. *Nat Commun* 2018, 9:1069.
- 15 21. Møller HD: Circle-Seq: Isolation and Sequencing of Chromosome-Derived Circular DNA Elements in  
16 Cells. *Methods Mol Biol* 2020, 2119:165-181.
- 17 22. deCarvalho AC, Kim H, Poisson LM, Winn ME, Mueller C, Cherba D, Koeman J, Seth S, Protopopov A,  
18 Felicella M, et al: Discordant inheritance of chromosomal and extrachromosomal DNA elements  
19 contributes to dynamic disease evolution in glioblastoma. *Nat Genet* 2018, 50:708-717.
- 20 23. Mehta D, Cornet L, Hirsch-Hoffmann M, Zaidi SS-e-A, Vanderschuren H: Full-length sequencing of  
21 circular DNA viruses and extrachromosomal circular DNA using CIDER-Seq. *Nature Protocols* 2020,



- 1 15:1673-1689.
- 2 24. Buenrostro JD, Wu B, Chang HY, Greenleaf WJ: ATAC-seq: A Method for Assaying Chromatin  
3 Accessibility Genome-Wide. *Curr Protoc Mol Biol* 2015, 109:21.29.21-21.29.29.
- 4 25. Schones DE, Cui K, Cuddapah S, Roh TY, Barski A, Wang Z, Wei G, Zhao K: Dynamic regulation of  
5 nucleosome positioning in the human genome. *Cell* 2008, 132:887-898.
- 6 26. Lee I, Razaghi R, Gilpatrick T, Molnar M, Gershman A, Sadowski N, Sedlazeck FJ, Hansen KD,  
7 Simpson JT, Timp W: Simultaneous profiling of chromatin accessibility and methylation on human cell  
8 lines with nanopore sequencing. *Nature Methods* 2020, 17:1191-1199.
- 9 27. Shipony Z, Marinov GK, Swaffer MP, Sinnott-Armstrong NA, Skotheim JM, Kundaje A, Greenleaf WJ:  
10 Long-range single-molecule mapping of chromatin accessibility in eukaryotes. *Nat Methods* 2020,  
11 17:319-327.
- 12 28. Stergachis AB, Debo BM, Haugen E, Churchman LS, Stamatoyannopoulos JA: Single-molecule  
13 regulatory architectures captured by chromatin fiber sequencing. *Science* 2020, 368:1449-1454 %@  
14 0036-8075.
- 15 29. Gaubatz JW, Flores SC: Purification of eucaryotic extrachromosomal circular DNAs using exonuclease  
16 III. *Anal Biochem* 1990, 184:305-310.
- 17 30. Huddleston J, Chaisson MJP, Steinberg KM, Warren W, Hoekzema K, Gordon D, Graves-Lindsay TA,  
18 Munson KM, Kronenberg ZN, Vives L, et al: Discovery and genotyping of structural variation from  
19 long-read haploid genome sequence data. *Genome Res* 2017, 27:677-685.
- 20 31. He HH, Meyer CA, Chen MW, Jordan VC, Brown M, Liu XS: Differential DNase I hypersensitivity  
21 reveals factor-dependent chromatin dynamics. *Genome Res* 2012, 22:1015-1025.

- 1 32. Møller HD, Mohiyuddin M, Prada-Luengo I, Sailani MR, Halling JF, Plomgaard P, Maretty L, Hansen  
2 AJ, Snyder MP, Pilegaard H, et al: Circular DNA elements of chromosomal origin are common in  
3 healthy human somatic tissue. *Nature Communications* 2018, 9:1069.
- 4 33. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS,  
5 Eppig JT, et al: Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat*  
6 *Genet* 2000, 25:25-29.
- 7 34. Barutcu AR, Lajoie BR, McCord RP, Tye CE, Hong D, Messier TL, Browne G, van Wijnen AJ, Lian JB,  
8 Stein JL, et al: Chromatin interaction analysis reveals changes in small chromosome and telomere  
9 clustering between epithelial and breast cancer cells. *Genome Biol* 2015, 16:214.
- 10 35. Klemm SL, Shipony Z, Greenleaf WJ: Chromatin accessibility and the regulatory epigenome. *Nat Rev*  
11 *Genet* 2019, 20:207-220.
- 12 36. Conconi A, Widmer RM, Koller T, Sogo J: Two different chromatin structures coexist in ribosomal RNA  
13 genes throughout the cell cycle. *Cell* 1989, 57:753-761.
- 14 37. Moller HD, Parsons L, Jorgensen TS, Botstein D, Regenberg B: Extrachromosomal circular DNA is  
15 common in yeast. *Proc Natl Acad Sci U S A* 2015, 112:E3114-3122.
- 16 38. Smith CL, Cantor CR: 6 - Purification, Specific Fragmentation, and Separation of Large DNA  
17 Molecules. In *Recombinant DNA Methodology*. Edited by Wu R, Grossman L, Moldave K. San Diego:  
18 Academic Press; 1989: 139-157
- 19 39. Li H: Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018, 34:3094-3100.
- 20 40. Guo X, Chen F, Gao F, Li L, Liu K, You L, Hua C, Yang F, Liu W, Peng C: CNSA: a data repository for  
21 archiving omics data. *Database* 2020, 2020.

1 41. Chen FZ, You LJ, Yang F, et al. CNGBdb: China National GeneBank DataBase, *Hereditas*.

2 2020;42(08):799-809. doi:10.16288/j.ycz.20-080

3

4

5

6

7

8

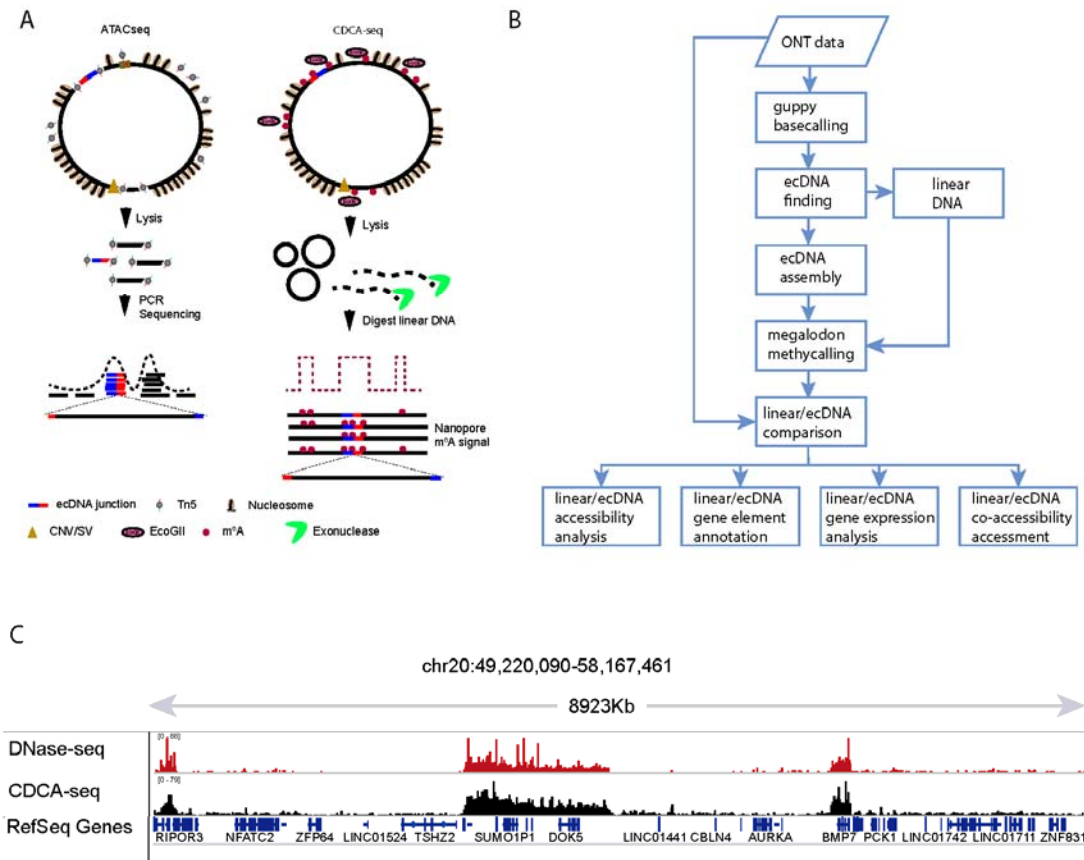
9

10

11

12

13



1

2 **Figure1. The CDCA-seq for profiling chromatin accessibility and nucleosome position**

3 **at ecDNAs.** (A) Intact chromatin is treated with m<sup>6</sup>A methyltransferase (EcoGII), which

4 preferentially methylate DNA bases in open chromatin region on ecDNAs and linear DNAs.

5 High molecular weight DNA is then isolated and subjected to exonuclease digestion to

6 remove partially linear DNAs. The rest DNAs are subjected to nanopore library construction

7 and nanopore sequencing. The data were aligned to genome to identify ecDNAs based on

8 head-to-tail pattern. The methylated bases are used to reconstruct nucleosomes in ecDNAs

9 and rest linear DNAs. In contrast, the ATAC-seq used the transposon to attack the open

10 chromatin. The tagmented short fragments are amplified and subjected to NGS. The short

1 reads are aligned with genome to identify ecDNAs based. The mapped reads are calling as  
2 peaks to represent the open chromatin region. (B) The bioinformatics pipeline of CDCA-seq.  
3 The signal data were processed through guppy basecalling to generate sequence. The  
4 sequences were aligned to genome to identify the linear DNA and ecDNAs. We assembled  
5 the ecDNA sequence reference. Based on the ecDNA and linear DNA reference, we used  
6 Megalodon to call the m<sup>6</sup>A sites base on ecDNA and linear DNAs. Then we perform the  
7 accessibility analysis, gene element annotation, gene expression analysis, co-accessibility  
8 assessment. (C) Large aggregate CDCA-seq signal enrichments match closely with  
9 DNase-seq accessibility peaks. (Chr20:49,220,090-58,167,461)

10

11

12

13

14

15

16

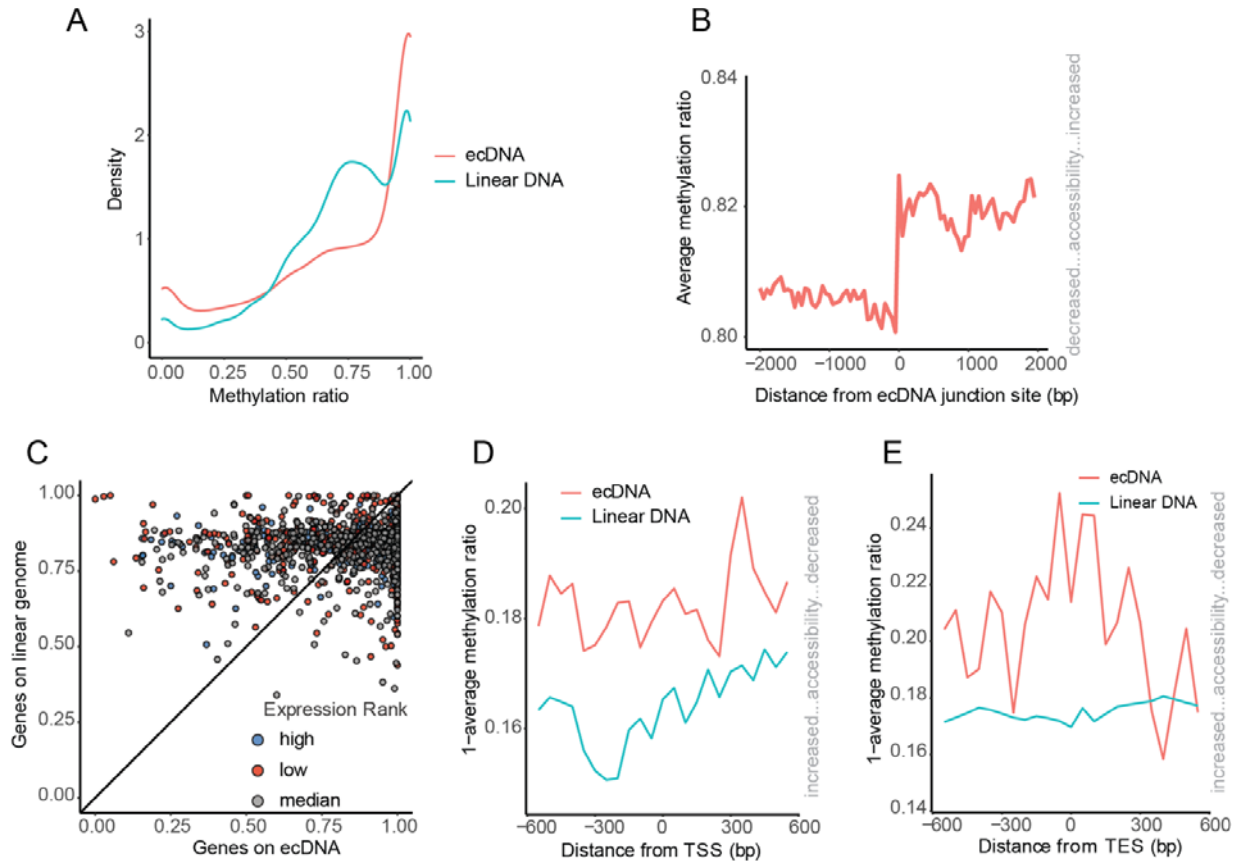
17

18

19

20

21



1

2 **Figure 2. The ecDNA and linear DNA have the different chromatin accessibility pattern.**

3 (A) The density distribution of the methylation ratio on ecDNAs and linear DNAs. (B) The

4 average chromatin status around ecDNA neighboring regions. The junction site and its right

5 neighboring regions demonstrate the more open chromatin. (C) The average methylations of

6 gene regions on ecDNA and linear DNAs (from TSS to TES). The genes were classified as

7 two groups: I. The genes on linear DNA have more open chromatin structure than ecDNA

8 carried genes; II. The ecDNA carried genes have more open chromatin structure than the

9 genes on linear DNA. (D) Average CDCA-seq profile around transcription start site on

10 ecDNAs and linear DNAs. (E) Average CDCA-seq profile around transcription end site on

1 ecDNAs and linear DNAs. (aggregated over 50-bp windows sliding every 5 bp; the  
2 sequencing depth is normalized for ecDNA and linear DNA;see Methods for details)

3

4

5

6

7

8

9

10

11

12

13

14

15

16

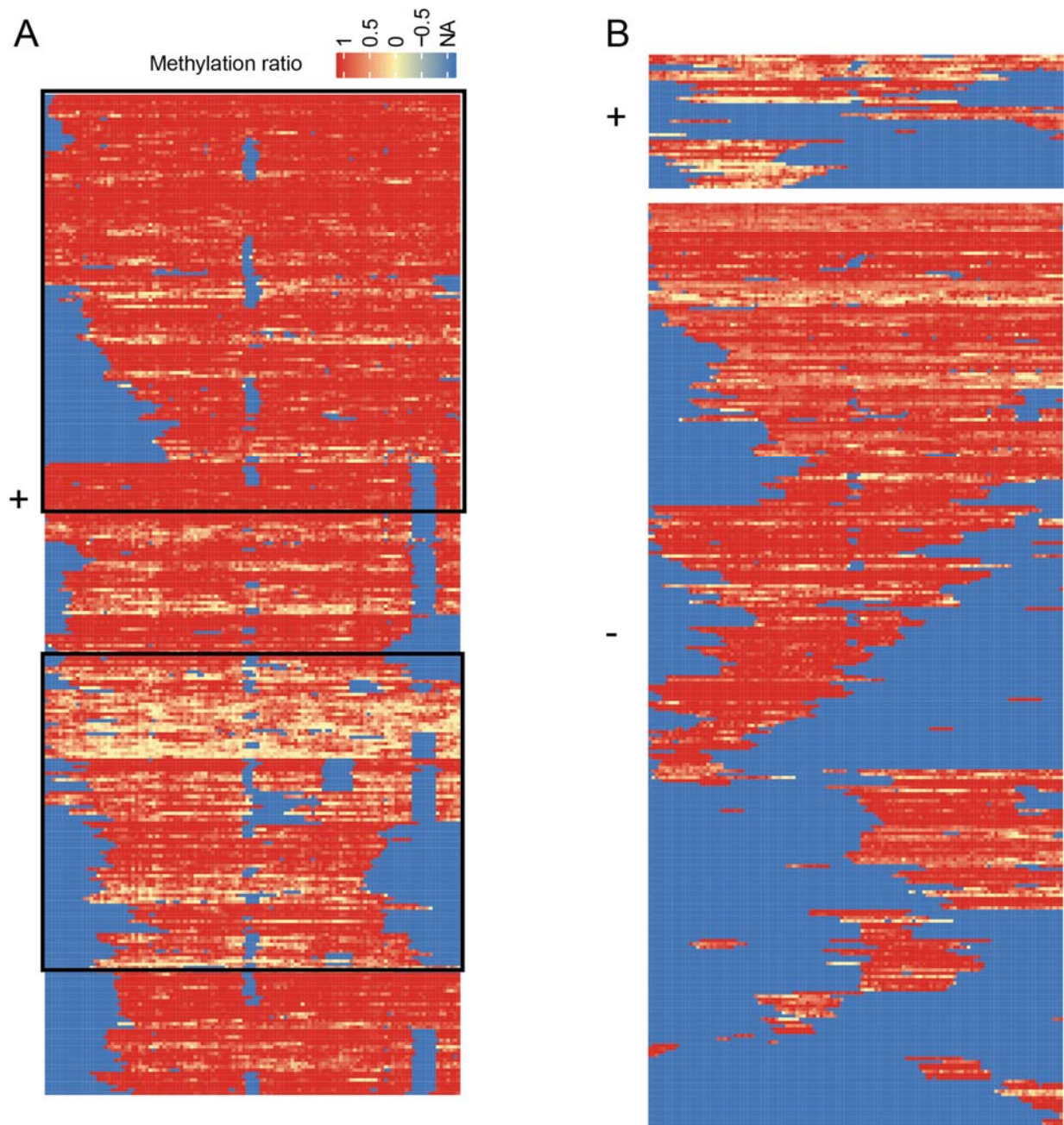
17

18

19

20

21



1

2 **Figure 3. CDCA-seq reveals the distribution of alternative chromatin states of**

3 **ecDNA arrays.** A. Shown are all reads covering the linear DNA region

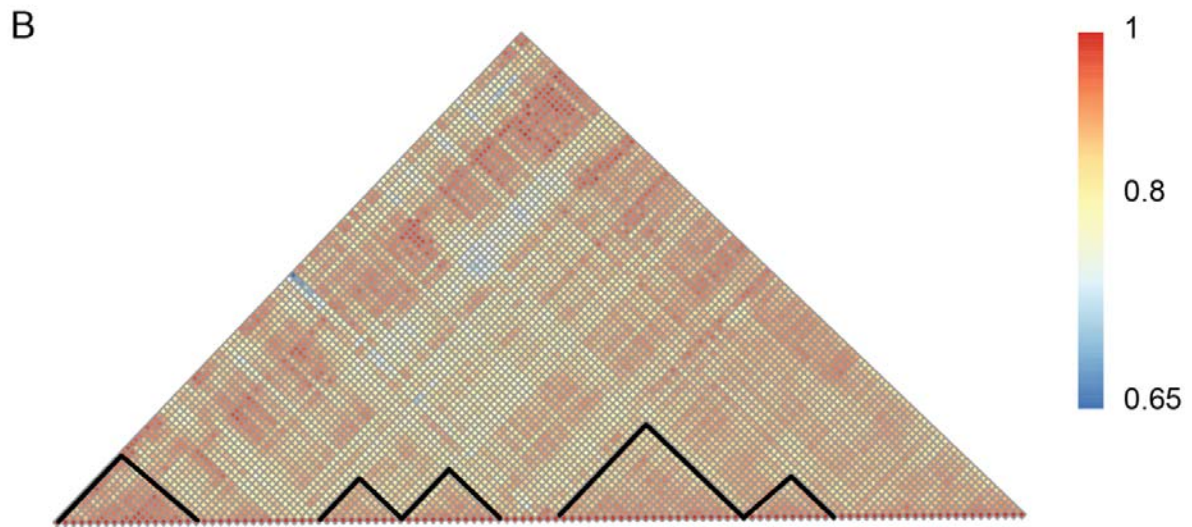
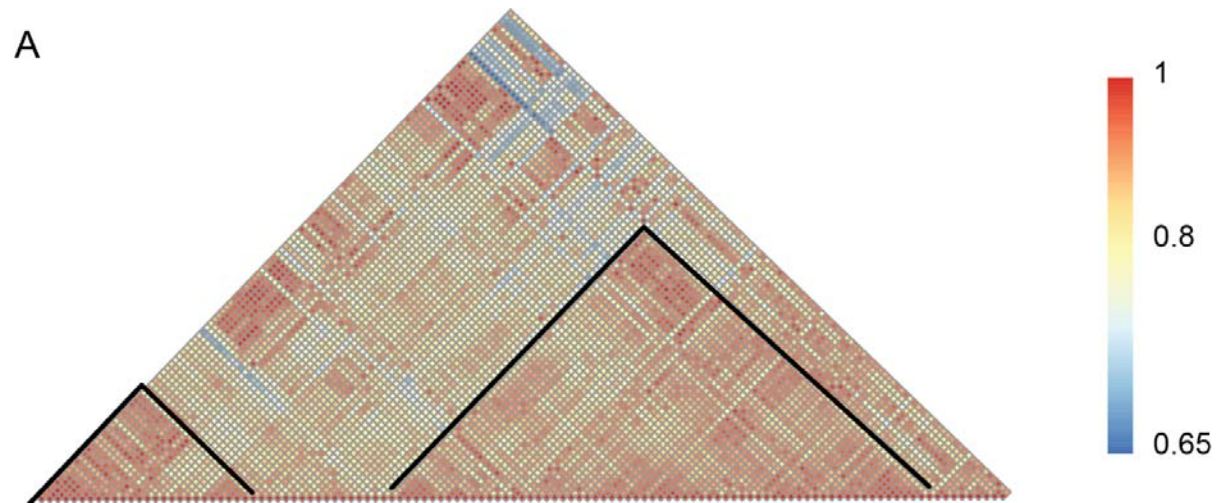
4 chr10:42383201~42389201. The box highlights the active and inactive chromatin. B.

5 Shown are all reads covering the ecDNA region chr10:42383201~42389201. The upper

6 panel indicates the positive strand, and the lower panel indicates the negative strand.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19



1  
2 **Figure 4. Chromatin co-accessibility profiles for the chr10:42383201~42389201**  
3 **show correlation and anticorrelation in ecDNA and linear DNA. (A) Chromatin**  
4 **co-accessibility profiles for the chr10:42383201~42389201 show correlation and**  
5 **anticorrelation on ecDNA. Red indicates the positive correlation and blue indicates the**  
6 **anticorrelation. (B) Chromatin co-accessibility profiles for the**  
7 **chr10:42383201~42389201 show correlation and anticorrelation on linear DNA.**