

1 **CAM: a quality control pipeline for MNase-seq data**

2

3 Sheng'en Hu ^{1¶}, Xiaolan Chen^{1¶}, Ji Liao¹, Yiqing Chen¹, Chengchen Zhao¹, Yong
4 Zhang^{1*}

5

6

7 ¹Translational Medical Center for Stem Cell Therapy & Institute for Regenerative Medicine,
8 Shanghai East Hospital, School of Life Science and Technology, Shanghai Key Laboratory of
9 Signaling and Disease Research, Tongji University, Shanghai 20092, China

10

11 ¶ These two authors contributed equally to this work.

12 *Corresponding author.

13 Email: yzhang@tongji.edu.cn

14

15

16

17 **Abstract**

18 Nucleosome organization affects the accessibility of cis-elements to trans-acting
19 factors. Micrococcal nuclease digestion followed by high-throughput sequencing
20 (MNase-seq) is the most popular technology used to profile nucleosome organization
21 on a genome-wide scale. Evaluating the data quality of MNase-seq data remains
22 challenging, especially in mammalian. There is a strong need for a convenient and
23 comprehensive approach to obtain dedicated quality control (QC) for MNase-seq data
24 analysis. Here we developed CAM, which is a comprehensive QC pipeline for
25 MNase-seq data. The CAM pipeline provides multiple informative QC measurements
26 and nucleosome organization profiles on different potentially functional regions for
27 given MNase-seq data. CAM also includes 268 historical MNase-seq datasets from
28 human and mouse as a reference atlas for unbiased assessment. CAM is freely
29 available at: <http://www.tongji.edu.cn/~zhanglab/CAM>

30

31 **Keywords:** MNase-seq, Quality control pipeline, Transcriptional regulation

32

33 **Introduction**

34 Nucleosome organization (i.e., the relative location of a nucleosome on the DNA)
35 affects the transcriptional activity by influencing the access of DNA-binding proteins
36 to the genome and the elongation of RNA polymerase II [1, 2]. Recently, nucleosome
37 organizations in a variety of species and cell types have been profiled using
38 micrococcal nuclease digestion followed by high-throughput sequencing (MNase-seq)
39 [3]. Although MNase-seq technology has been widely used and many computational
40 tools have been developed for MNase-seq data [4], the quality evaluation remains
41 challenging, especially for data from mammalian genomes, due to two major
42 difficulties. First, different experimental designs (e.g., sequencing coverage and the
43 MNase concentration) may result in distinct nucleosome organization features in
44 some genomic loci (e.g., fragile nucleosomes at promoters [5]). Second, in contrast to

45 chromatin immunoprecipitation sequencing (ChIP-seq), DNase-seq and methylated
46 DNA immunoprecipitation sequencing (MeDIP-seq) data, the MNase-seq data signals
47 are not enriched in any specific genomic loci, resulting in difficulties in focusing on
48 target regions for downstream analysis. Many software tools were designed to detect
49 well-positioned nucleosomes, but seldom took care of MNase specific quality control
50 which is the basis for detecting nucleosome organization correctly and precisely. Here,
51 we present CAM, an integrated quality control (QC) for MNase-seq data. CAM
52 provides multiple key measurements that enable users to evaluate the data quality
53 using scores from 268 historical MNase-seq datasets in human and mouse as a
54 reference atlas. In addition, CAM provides nucleosome organization information
55 based on potentially functionally related genomic regions for use in the targeted
56 downstream analysis.

57

58 **Results and conclusion**

59 **Overview of CAM**

60

61 The CAM pipeline initiates from the data pre-processing steps, including reads
62 mapping (optional), high-quality reads filtering (optional) and nucleosome
63 organization profile generation (Fig 1). After the pre-processing steps, CAM provides
64 multiple QC measurements to allow users to evaluate the data quality as follows: 1)
65 sequencing coverage, 2) AA/TT/AT dinucleotide frequency, 3) nucleosomal DNA
66 length, 4) existence of nucleosome free regions (NFR) at promoters, 5) well-
67 positioned nucleosomes at the downstream promoters, 6) well-positioned
68 nucleosomes at custom defined potential cis-regulatory regions, 7) enrichment of
69 well-positioned nucleosome arrays in DNase hypersensitive sites (DHS) (S1 Table).
70 We compiled 268 MNase-seq datasets from human and mouse as a historical QC
71 reference atlas and made an unbiased judgment for the measurements by generating
72 several QC scores (Method section, S2 Table). In addition, CAM provides related
73 analysis results including 1) nucleosome profiles at each of promoter and custom
74 defined region, 2) detection of well-positioned nucleosome arrays, 3) gene level
75 annotation of the well-positioned nucleosome arrays.

76

77 **Fig. 1. Flowchart illustrating the CAM with default parameters.** The workflow of
78 CAM includes three components: data pre-processing (black box), analysis (grey box)
79 and QC (white box). The QC metrics for all historical MNase-seq data were
80 precompiled for unbiased judgment in the QC component.

81

82 Many software tools have been developed for MNase-seq data analysis. However,
83 most of them focused on detecting positioned nucleosomes and did not take care of
84 the quality of MNase-seq data, which has great importance on nucleosome detection.
85 Additionally, compared with the concept of positioned nucleosome, array of well-
86 positioned nucleosome provides a more unbiased description of nucleosome
87 positioning detection for MNase-seq data [1, 6-8] (discussed in the following section)
88 (S3 Table). CAM has three advantages over existing state-of-the-art methods: 1)
89 CAM provides QC components specific to MNase-seq. 2) Like most of existing
90 software tools, CAM also provides a peak calling function, however the peak calling
91 function is based on a more unbiased concept: well-positioned nucleosome array. 3)
92 CAM is implemented as a systematic pipeline, which takes raw MNase-seq reads or
93 aligned reads as input, and outputs QC measurements together with a series of
94 analysis results (discussed in the following section).

95

96 **Historical QC reference atlas from published MNase-seq** 97 **data**

98

99 To provide an unbiased judgment of QC measurements, we collected 268 MNase-seq
100 samples from human and mouse and pre-processed with all of our QC components.
101 Sequencing coverage, rotational score, estimated fragment length and nucleosome
102 profiles at promoters were collected as historical QC reference to measure the quality
103 of input MNase-seq data (S2 Table).

104

105 **Sequencing coverage measures the resolution of nucleosome** 106 **positioning**

107

108 Sequencing coverage provides a direct measurement of the resolution of two features
109 of nucleosome organization (i.e., occupancy and positioning) [9]. We generated
110 simulated regions with perfect positioning and no positioning, and compared the
111 difference of nucleosome positioning scores with different sequencing coverage to

112 show the influence of sequencing coverage on the resolution of nucleosome
113 positioning detection (Method section). The perfect positioning region and no
114 positioning region showed similar positioning scores with low sequencing coverage
115 (almost the same between them in 1-fold coverage) and the nucleosome positioning
116 score of the perfect positioning region significantly increased compared with the no
117 positioning region with higher sequencing coverage (Fig 2A). In the CAM pipeline,
118 we calculated sequencing coverage of the input MNase-seq data and compared it with
119 the distribution of historical data (Fig 2B) to reflect the ability of MNase-seq data to
120 detect nucleosome positioning and occupancy.

121

122 **Fig. 2. Average nucleosome coverage. (A)** Regions with higher sequencing coverage
123 exhibit higher resolution for the detection of nucleosome positioning. The barplot
124 shows the nucleosome positioning score in a simulated perfect positioning region and
125 no positioning region. For low sequencing coverage, the positioning score is almost
126 the same for the perfect and no positioning regions, and the difference increases when
127 the coverage increases. **(B)** Distribution of sequencing coverage of all historical data.

128

129 **AA/TT/AT periodicity measures the nucleosome rotational** 130 **positioning**

131

132 The 10-base AA/TT/AT di-nucleotide periodicity in the nucleosomal DNA provides a
133 measurement of nucleosome rotational positioning, which is influenced by the DNA
134 sequence. Studies show that AA/TT/AT di-nucleotide displays a 10 bp periodicity
135 throughout nucleosomal DNA sequence. The di-nucleotide favors DNA bending in a
136 specific direction and expands the major groove [10]. Thus the 10 bp periodicity
137 should be observed in the AA/TT/AT di-nucleotide frequency of MNase-seq reads if
138 the MNase-seq reads reflect nucleosomal DNA accurately. We defined a “rotational
139 score” to measure the 10 bp AA/TT/AT periodicity (Method section) and compared
140 with the distribution of historical data (Fig 3A). A judgment of “rotational score”
141 (Pass or Fail) was assigned for the input sample to measure the rotational positioning
142 (Fig 3B).

143

144 **Fig. 3. AA/TT/AT di-nucleotide periodicity. (A)** Distribution of the rotational

145 scores of all historical data. A rotational score < 0.08 was determined as a “Fail” (blue)
146 in this measurement. **(B)** Examples of “Pass” and “Fail” MNase-seq data in this
147 measurement.

148

149 **The nucleosomal DNA length distribution reflects the** 150 **MNase concentration**

151 The nucleosomal DNA length distribution (referring to the fragment length or MNase
152 library size) is closely related and thus reflects the MNase concentration. MNase
153 concentration is negatively correlated with the fragment length estimated from the
154 MNase-seq data: Samples with higher MNase concentration exhibited shorter
155 nucleosomal DNA length, whereas lower concentration samples exhibited longer
156 nucleosomal DNA length (Fig 4A). CAM estimates fragment length of nucleosomal
157 DNA from the input data and compared with the distribution of historical data to
158 reflect the MNase concentration (Method section, Figure 4B). Although different
159 MNase concentration may result in different nucleosome organization features in
160 some genomic loci (e.g. fragile nucleosomes at promoters) [5], we expect the
161 fragment length to be close to the length of nucleosomal DNA (147 bp), indicating
162 that the MNase-seq reads accurately represent the location of mono-nucleosome
163 (through a judgment of Pass or Fail, Fig 4B, C). A longer or shorter fragment length
164 indicates partial or over digestion of MNase, respectively.

165

166 **Fig. 4. Nucleosomal DNA length distribution.** **(A)** The MNase digestion level
167 (concentration) is related to the nucleosomal DNA length: samples with a higher
168 MNase concentration exhibit shorter nucleosomal DNA length, whereas lower
169 concentration samples exhibit longer nucleosomal DNA length. The MNase-seq data
170 used for the comparison were obtained from a previous study in mouse ESCs
171 (GSM2083105, GSM2083106, GSM1083107, GSM1083108). **(B)** The distribution of
172 the nucleosome length from all of the historical data. Nucleosome length < 140 or $>$
173 155 was determined as “Fail” (blue) in this measurement. **(C)** Examples of MNase-
174 seq data with “Pass” and “Fail” in this measurement. The vertical line labels 147 bp.

175

176 **Nucleosome profiles on potentially functional regions reflect**

177 **the ability for detecting well-positioned nucleosomes**

178 Well-positioned nucleosome, an effective marker of transcriptional regulatory regions,
179 is formed by nucleosomes consistently positioned in a cell population. Thus, MNase-seq
180 reads can be observed to position consistently in certain functional regions to
181 form well-positioned nucleosomes. Promoter regions are regarded as the most
182 important functional regions for transcriptional regulation. According to the barrier
183 nucleosome model [1, 2], nucleosome free regions (NFR) and the successive well-
184 positioned nucleosomes are supposed to be observed around TSS. CAM generates
185 nucleosome profiles on promoter regions and displays with an aggregate plot and a
186 heatmap (Fig 5A). Additionally, CAM calculates two QC scores according to the
187 nucleosome profiles on promoter regions, and assigns judgments for the input sample
188 to measure the nucleosome positioning pattern at promoter regions, based on the
189 comparison with the distribution of the historical data (Method section). The two QC
190 scores consist of 1) a promoter NFR score to check the existence of NFR at promoters
191 (Fig 6A, B) and 2) a promoter positioning score to measure the successive well-
192 positioned nucleosomes at the downstream promoters (Fig 6C, D).

193

194 **Fig. 5. Nucleosome profiles on promoter regions and custom regions.** Nucleosome
195 profiles were generated and plotted as aggregate plots and heatmaps on (A) promoters
196 and (B) custom regions. By default, the promoter regions were defined as -1 kb to +2
197 kb from the TSS of all refseq genes, and the custom regions were extended to +/-1 kb
198 (by default) from the center of the regions. Both the promoters and custom regions
199 were sorted by the MNase-seq read counts within the regions. An MNase-seq data
200 from a human lymphoblastoid cell line (GSM907784) was selected to plot the
201 nucleosome profiles on both promoters and custom regions.

202

203 **Fig. 6. Distribution of promoter NFR score and promoter positioning score.** (A)
204 Distribution of the promoter NFR score of all historical data. A promoter NFR score
205 < 0.4 was determined as a “Fail” (blue) in this measurement. (B) Examples of “Pass”
206 and “Fail” MNase-seq data in promoter NFR score. (C) Distribution of the promoter
207 positioning scores of all historical data. A promoter positioning score > 0.4 was
208 determined as a “Fail” (blue) in this measurement. (D) Examples of “Pass” and “Fail”
209 MNase-seq data in promoter positioning score.

210

211 In addition, users can provide custom regions for CAM to profile nucleosome
212 organization at the same time. For example, users can provide motif sites or binding
213 sites (defined by ChIP-Seq peaks) of certain transcription factors (CTCF motif sites as
214 example, Fig 5B) as custom regions, and check whether the successive well-
215 positioned nucleosomes are also observed at these regions as an additional QC
216 measurement.

217 In our previous work, we developed an effective method to detect cis-regulatory
218 regions with MNase-seq data, which we called the detection of well-positioned
219 nucleosome arrays [8]. A well-positioned nucleosome array is a broader region with
220 successive well-positioned nucleosomes, often emanate from a nucleosome-depleted
221 region, such as transcription factor binding site [1, 6, 7], and is very unlikely to be
222 generated at random.

223 CAM adopts the method to detect nucleosome arrays across whole genome (Method
224 section, Fig 7A displays an example of detected well-positioned nucleosome arrays)
225 [8]. The nucleosome arrays define a list of potential target regions for the following
226 analysis, which is similar as ChIP-Seq peaks, and solve the problem that MNase-seq
227 data do not show enrichment in any regions. We showed that these nucleosome arrays
228 were enriched in potentially functional regions such as the downstream promoters and
229 the union DHS sites (Method section, Fig 7B), which can be explained by the barrier
230 nucleosome model [11]. Thus we calculated the fold enrichment of the well-
231 positioned nucleosome arrays on the union DHS sites compared with random
232 background (Method section) as a QC measurement to assess the ability of the
233 MNase-seq data to detect cis-regulatory elements with well-positioned nucleosome
234 arrays. Samples with fold enrichment less than 2 are regarded as “Fail” in this
235 measurement, indicating the well-positioned nucleosome arrays are more likely to be
236 caused by random rather than the barrier nucleosome model. The genomic coordinates
237 together with the nucleosome profile values are reported for each region with a well-
238 positioned nucleosome array. CAM also provides a gene level annotation of the well-
239 positioned nucleosome arrays for downstream analysis (Method section).

240

241 **Fig. 7. Well-positioned nucleosome array.** (A) An example of detected well-
242 positioned nucleosome arrays. An MNase-seq data from a human lymphoblastoid cell
243 line (GSM907784) was selected to plot the nucleosome profile for the example region.

244 **(B)** Well-positioned nucleosome arrays are enriched in both the downstream
245 promoters and the union DHS sites. The enrichment fold (y-axis) represents the fold
246 change of the observed proportion of nucleosome arrays on the downstream
247 promoters (or the union DHS sites) to the expected proportion (Method section). The
248 nucleosome arrays were detected using MNase-seq data from human lymphoblastoid
249 cell line (GSM907784).

250

251

252 **Computational cost and standard output of CAM**

253 We applied CAM to published MNase-seq data from a human lymphoblastoid cell
254 line (GSM907784) and obtained multifaceted and detailed QC reports (S1 File)
255 together with a series of analysis results (S4 Table). Running time the CAM pipeline
256 is also provided (S5 Table).

257

258 **Conclusion**

259 In summary, CAM is specifically designed for QC of MNase-seq data. The QC
260 components measure the quality of MNase-seq data in different aspects (S1 Table).
261 The program uses standard format input files via simple commands, reports
262 informative QC measurements (specific for MNase-seq data) to assist in evaluating
263 the data quality (using historical MNase-seq data as a reference atlas), generates
264 nucleosome organization profiles based on promoters and custom defined regions,
265 and detects regions with well-positioned nucleosome arrays.

266

267 **Materials and Methods**

268 **Implementation and webpages of CAM**

269 CAM was implemented using Python and R. Users need python (version = 2.7) and R
270 (version >= 2.14.1) installed on linux or MacOS environment to make sure the
271 successful process of CAM. CAM was distributed under the GNU General Public
272 License. The source code and detailed tutorial of CAM is available on our webpages:

273 <http://www.tongji.edu.cn/~zhanglab/CAM>.

274

275 **Data pre-processing**

276

277 In the alignment process, Bowtie (Langmead et al., 2009) was used to align the
278 MNase-seq reads with the -m 1 parameter. Users could adjust two other parameters
279 (“-X” for maximum insert size in paired end data and “-3” for trimming bases from
280 3’end of reads) in the configure file. To maintain high-quality alignment results, we
281 removed reads with a sequencing quality less than 30 (The default cutoff of MAPQ
282 filtering is 30). The alignment step was skipped when the aligned reads (SAM/BED
283 format) input was used. CAM was designed specifically for MNase-seq data from
284 human and mouse and only supports (by default) the hg38, hg19, mm10, and mm9
285 genome versions. Users can add custom genome versions according to the
286 instructions in the CAM manual.

287 To generate nucleosome profiles, the sequencing reads were transformed into
288 nucleosome reads as follows: for single end sequencing data, the reads were extended
289 to 147 bp in the 3’ direction; for paired end data, the paired end fragments were
290 extended to 73 bp in both the 5’ and 3’ directions from the fragment centers. Then,
291 the middle 73 bp centered on the extended fragments were compiled as the
292 nucleosome profile.

293

294

295 **Calculation of sequencing coverage**

296

297 The sequencing coverage (fold) was defined as $(\text{Number of reads} \times 194 \text{ bp}) /$
298 $(\text{Effective genome size})$. The “number of reads” was defined as the number of
299 mappable reads after MAPQ filtering for single end data and the number of fragments
300 for paired end data. Additionally, “194 bp” represented the total length of the
301 nucleosomal DNA and linker DNA, which was estimated from 268 historical MNase-
302 seq datasets. The “effective genome size” was defined as 2.7 billion nucleotides
303 $(2.7e9)$ for human and 1.87 billion nucleotides $(1.87e9)$ for mouse. The effective
304 genome size used here was smaller than the original genome size due to the repetitive
305 features on the chromosomes.

306

307 **Simulation of the perfect positioning and no positioning**

308 **regions**

309

310 We compared the nucleosome positioning scores on a simulated perfect positioning
311 region with a no positioning region with different sequencing coverage to demonstrate
312 the influence of the sequencing coverage on the resolution of nucleosome positioning.
313 A 1200-bp region was prepared for all simulations. Five “potential nucleosome
314 centers” were marked with a distance between adjacent centers equal to 194 bp
315 (nucleosomal + linker DNA length, estimated with the historical data). Simulated
316 reads (1 bp read to reflect the nucleosome center) were assigned evenly to each of the
317 5 “potential nucleosome centers” in the perfect positioning region with a 5 bp
318 perturbation. For the no positioning region, simulated reads were assigned randomly
319 on the 1200 bp region with equal probability and the same perturbation as the perfect
320 positioning region (5 bp). The nucleosome positioning scores for the different
321 sequencing coverage in the two types of regions were calculated using the method
322 described below (“Detect well-positioned nucleosome arrays” section).

323

324 **Measurement of AA/TT/AT di-nucleotide frequency and** 325 **definition of rotational score**

326

327 We sampled down all mappable reads to 10 million and extended each read to 147 bp
328 in the 3' direction. Then, the aggregate AA/TT/AT di-nucleotide frequency was
329 calculated across 4 to 143 bp of the extended reads. We conducted a Fourier transform
330 on the aggregate frequency and used the energy of 10-bp periodicity (defined as the
331 rotational score) to demonstrate the extent to which the MNase-seq reads reflect the
332 nucleosome organization. Samples with rotational scores greater than 0.08 were
333 defined as “Pass” in this measurement, whereas the other samples were defined as
334 “Fail”. The cutoff 0.08 was determined from the distribution of rotational scores
335 among all historical MNase-seq data.

336

337 **Calculation of nucleosomal DNA length distribution**

338

339 For the paired end samples, the fragment length distribution from all mappable
340 fragments was used to directly infer the nucleosomal DNA length distribution. For the
341 single end samples, we calculated a start-to-end distance to estimate the nucleosomal
342 DNA length as follows: mappable reads were sampled down to 10 million; then, we

343 calculated the distribution of the distance from the 5' end of each plus strand read to
344 all 5' ends of the minus strand reads (the start-to-end distance) within 250 bp
345 downstream (1kb downstream in figures for visualization). Duplicate reads were
346 discarded in this calculation. After the distribution of the start-to-end distance was
347 generated, the length with the highest frequency was defined as the estimated
348 nucleosomal DNA length of the MNase-seq data (for both the paired end and single
349 end data). Samples with estimated nucleosomal DNA lengths in the range of 140 to
350 155 bp were defined as "Pass", whereas the other samples were defined as "Fail". The
351 cutoff was determined from the distribution of the nucleosomal DNA lengths among
352 all historical MNase-seq data.

353

354

355 **Nucleosome organizations on the promoters and custom** 356 **regions**

357

358 The nucleosome profile was generated as described in the "Data pre-processing"
359 section. The nucleosome signal from 1 kb upstream to 2 kb downstream (default, user
360 can change the range with certain parameters in the configure file; see Manual for
361 details) from TSS was plotted as a heatmap and aggregate curve with a 10-bp
362 resolution. For custom regions, the nucleosome signal was plotted within +/- 1kb
363 (default) from the center of the region. The signals from the minus strand transcripts
364 and the minus strand custom regions were reversed in both the heatmap and the
365 aggregate curve.

366

367 **Calculation of promoter NFR score and promoter** 368 **positioning score**

369

370 To calculate the promoter NFR score, we first generated the average MNase-seq
371 signal profile from 1kb upstream to 2kb downstream from all TSS. The aggregate
372 signal was then kernel smoothed to remove signal noise. Promoter NFR score was
373 calculated by the smoothed aggregate signal of the +1 nucleosome and the -1
374 nucleosome subtracting the signal of nucleosome free region. Finally, the promoter
375 NFR score was normalized by the difference between the maximum signal and the
376 minimum signal within the 3kb promoter regions. We regarded aggregate signal of the
377 first downstream local maximum as the signal of +1 nucleosome, the first upstream
378 local maximum as -1 nucleosome and the first upstream local minimum as NFR.

379 The promoter positioning score was defined as the coefficient of variance (CV) of the
380 distance between +1, +2, +3 and +4 nucleosomes. The position of nucleosomes was
381 defined as the local maximum positions.

382

383

384 **Detection of well-positioned nucleosome arrays**

385

386 To detect well-positioned nucleosome arrays, first, the mappable read were extended
387 to 147 bp in the 3' direction; then, the centers of the extended reads were compiled to
388 generate the nucleosome center profile (for the paired end data the center of each
389 fragment was compiled directly). Next, Gaussian smoothing (window size = +/-73 bp;
390 standard deviation = 30) was performed on the nucleosome center profile, and the
391 absolute difference between the adjacent bps was calculated as the modified profile.
392 Then, the local maximum within +/-73 bp was selected. All adjacent local maxima
393 were connected to create a “nuc-array” curve and a signal that was defined as the
394 “nuc-array” value. The “nuc-array” value was transformed to a fold enrichment value
395 over the background (defined as the average “nuc-array” value across the genome).
396 The well-positioned nucleosome array was defined as segments with lengths greater
397 than 600 bp and fold enrichment greater than 2 (default, the cutoff of length and fold
398 enrichment can be changed by the users). We generated a 5 columns bed file with the
399 above method (named as “outputname_Nucleosome_Array.bed”). Each line of the
400 bed file represented a well-positioned nucleosome array. The 5th column of the bed
401 file is the “nuc-array” value. The method was also described in a previous work [8].
402 The positioning score for the simulated regions (described above) was calculated by
403 the average “nuc-array” value across the whole simulated regions.

404

405

406 **Enrichment of well-positioned nucleosome arrays in** 407 **downstream promoter regions and union DHS sites**

408

409 The enrichment fold on the promoters of the nucleosome array was defined as the fold
410 change of observed and expected percentage of the well-positioned nucleosome
411 arrays on promoters. The expected percentage was equal to the percentage of the
412 promoter length compared with the effective genome size (mentioned above). Then,
413 the promoter regions were defined as 3 kb downstream from the TSS of refseq genes.
414 Similar enrichment fold was also performed on the Union DHS sites, which were

415 generated by merging the narrow peak of all DNase-seq data from the human and
416 mouse (separately) from the ENCODE and Roadmap Epigenomics project as
417 previously described [12]. Samples with enrichment fold greater than 2 were defined
418 as “Pass” in this measurement, otherwise they were “Fail”.

419

420 **Gene level annotation of well-positioned nucleosome arrays**

421 Well-positioned nucleosome arrays are annotated to genes by overlapping with the
422 promoters of genes, which are defined as 3 kb upstream and downstream from TSS.
423 The length of overlapped nucleosome arrays is marked as a feature of the certain gene
424 (“promoter nuc-array length”). CAM provides an additional analysis result in the
425 output folder, which is named as “outputname_geneLevel_nucarrayAnnotation.bed”.
426 The additional analysis result is a bed file with 6 columns, including “chromosome
427 name”, “promoter start site”, “promoter end site”, “refseq ID”, “promoter nuc-array
428 length” and “strand”. Each row of the bed file represents a refseq transcript. The 5th
429 column (“promoter nuc-array length”) represents the length of the well-positioned
430 nucleosome array which overlaps with the certain gene promoter (-1 for no
431 nucleosome array overlapped). With the help of this analysis result from different
432 samples, users can easily compare the array length of same genes in different
433 conditions.

434 **Acknowledgements**

435 We thank Kai Fu and Jiangxing Feng for their contributions in the early stage of this
436 project.

437

438 **Funding**

439 This work was supported by National Natural Science Foundation of China
440 (31571365, 31322031 and 31371288), National Key Research and Development
441 Program of China (2016YFA0100400), Specialized Research Fund for the Doctoral
442 Program of Higher Education (20130072110032), and Program of Shanghai
443 Academic Research Leader (17XD1403600).

444 *Conflict of Interest:* none declared.

445

446

447 **References**

448

- 449 1. Jiang C, Pugh BF. Nucleosome positioning and gene regulation: advances through genomics.
450 Nature reviews Genetics. 2009;10(3):161-72. doi: 10.1038/nrg2522. PubMed PMID: 19204718;
451 PubMed Central PMCID: PMC4860946.
- 452 2. Teves SS, Weber CM, Henikoff S. Transcribing through the nucleosome. Trends in biochemical
453 sciences. 2014;39(12):577-86. doi: 10.1016/j.tibs.2014.10.004. PubMed PMID: 25455758.
- 454 3. Hughes AL, Rando OJ. Mechanisms underlying nucleosome positioning in vivo. Annual review
455 of biophysics. 2014;43:41-63. doi: 10.1146/annurev-biophys-051013-023114. PubMed PMID:
456 24702039.
- 457 4. Teif VB. Nucleosome positioning: resources and tools online. Briefings in bioinformatics.
458 2016;17(5):745-57. doi: 10.1093/bib/bbv086. PubMed PMID: 26411474.
- 459 5. Xi Y, Yao J, Chen R, Li W, He X. Nucleosome fragility reveals novel functional states of
460 chromatin and poises genes for activation. Genome research. 2011;21(5):718-24. doi:
461 10.1101/gr.117101.110. PubMed PMID: 21363969; PubMed Central PMCID: PMC3083088.
- 462 6. Albert I, Mavrich TN, Tomsho LP, Qi J, Zanton SJ, Schuster SC, et al. Translational and
463 rotational settings of H2A.Z nucleosomes across the *Saccharomyces cerevisiae* genome. Nature.
464 2007;446(7135):572-6. doi: 10.1038/nature05632. PubMed PMID: 17392789.
- 465 7. Valouev A, Johnson SM, Boyd SD, Smith CL, Fire AZ, Sidow A. Determinants of nucleosome
466 organization in primary human cells. Nature. 2011;474(7352):516-20. doi: 10.1038/nature10002.
467 PubMed PMID: 21602827; PubMed Central PMCID: PMC3212987.
- 468 8. Zhang Y, Vastenhouw NL, Feng J, Fu K, Wang C, Ge Y, et al. Canonical nucleosome organization
469 at promoters forms during genome activation. Genome research. 2014;24(2):260-6. doi:
470 10.1101/gr.157750.113. PubMed PMID: 24285721; PubMed Central PMCID: PMC3912416.
- 471 9. Struhl K, Segal E. Determinants of nucleosome positioning. Nature structural & molecular
472 biology. 2013;20(3):267-73. doi: 10.1038/nsmb.2506. PubMed PMID: 23463311; PubMed Central
473 PMCID: PMC3740156.
- 474 10. Satchwell SC, Drew HR, Travers AA. Sequence periodicities in chicken nucleosome core DNA.
475 Journal of molecular biology. 1986;191(4):659-75. PubMed PMID: 3806678.
- 476 11. Mavrich TN, Ioshikhes IP, Venters BJ, Jiang C, Tomsho LP, Qi J, et al. A barrier nucleosome
477 model for statistical positioning of nucleosomes throughout the yeast genome. Genome research.
478 2008;18(7):1073-83. doi: 10.1101/gr.078261.108. PubMed PMID: 18550805; PubMed Central PMCID:
479 PMC2493396.
- 480 12. Zang C, Wang T, Deng K, Li B, Hu S, Qin Q, et al. High-dimensional genomic data bias
481 correction and data integration using MANCIE. Nature communications. 2016;7:11305. doi:
482 10.1038/ncomms11305. PubMed PMID: 27072482; PubMed Central PMCID: PMC4833864.
- 483

484 **Supporting Information**

485 **S1 File. CAM QC reports with default parameter for MNase-seq data from**
486 **human lymphoblastoid cell line (GSM907784).**

487 **S1 Table. Summary of QC measurements in CAM.**

488 **S2 Table. Meta data, QC scores and accession ID for 268 samples of historical**

489 **MNase-seq data.** Sequencing coverage, rotational scores, estimated fragment length,

490 promoter NFR scores and promoter positioning scores. The scores and judgments for

491 historical MNase-seq data were also attached.

492 **S3 Table. Compare of function between CAM and other nucleosome analysis**

493 **software.** We compare the major function of CAM to existing state-of-the-art

494 methods. CAM has three advantages: 1) Quality control component specific to

495 MNase-seq, 2) Well-positioned nucleosome arrays detection and 3) Systematic

496 pipeline for users to input raw sequencing data and get all QC and analysis results.

497 **S4 Table. List of standard output from CAM.**

498 **S5 Table. Running time of CAM.** An MNase-seq data from human lymphoblastoid

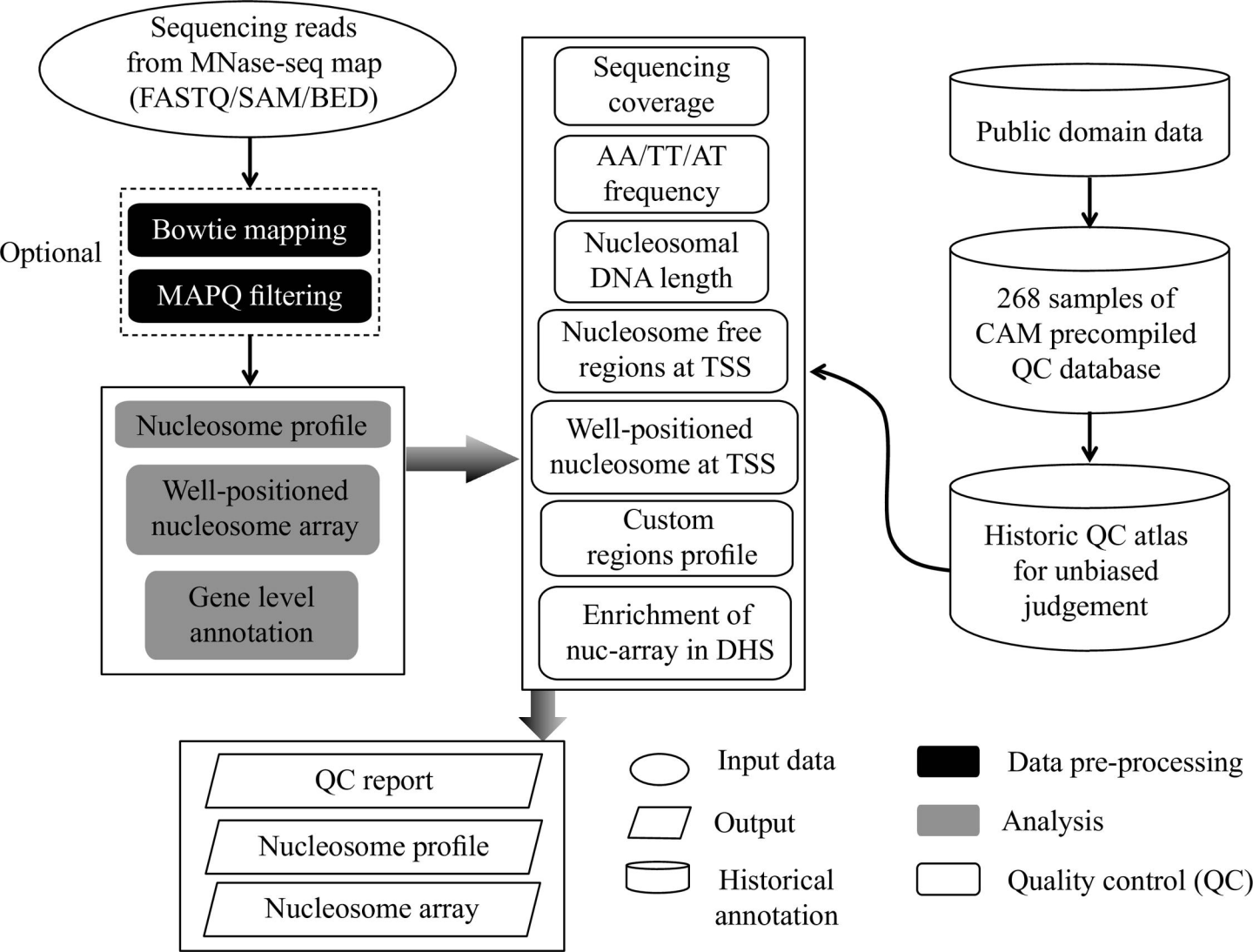
499 cell line (GSM907784, totally 546,924,994 reads) was used to evaluate the runtime of

500 CAM. Alignment process was excluded from this calculation. The percentage running

501 time for each component was calculated by using single CPU (Intel® Xeon® CPU

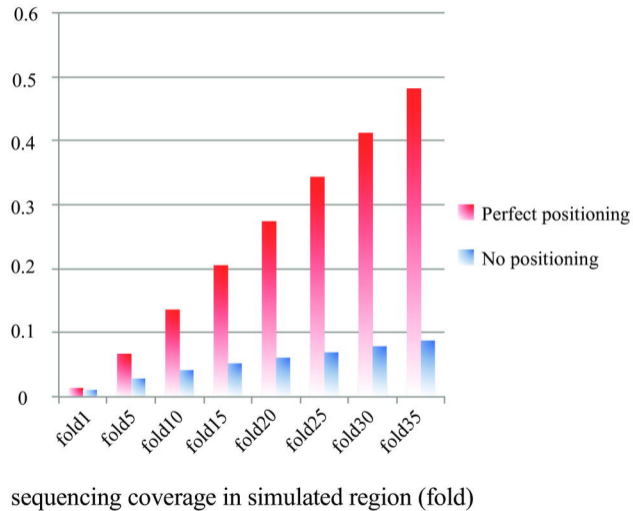
502 E5-2640 v2 @ 2.00 GHz).

503

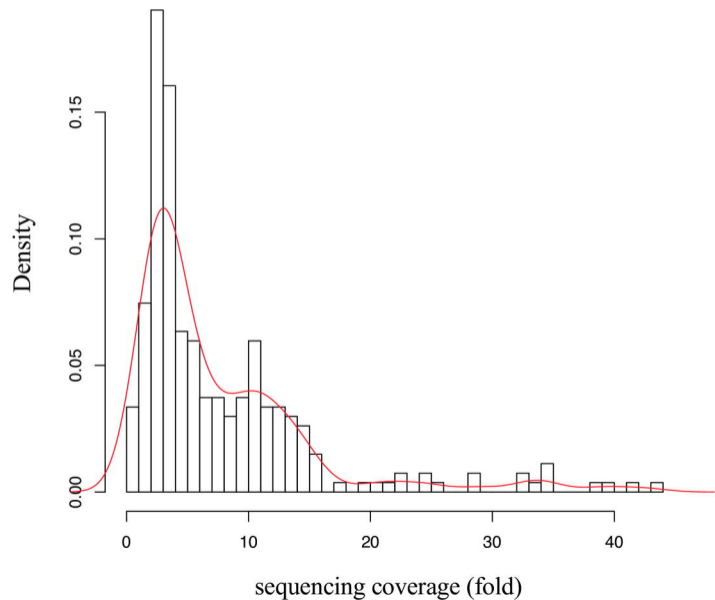


A

nucleosome positioning score on simulated
perfect and no position region
under different sequencing coverage

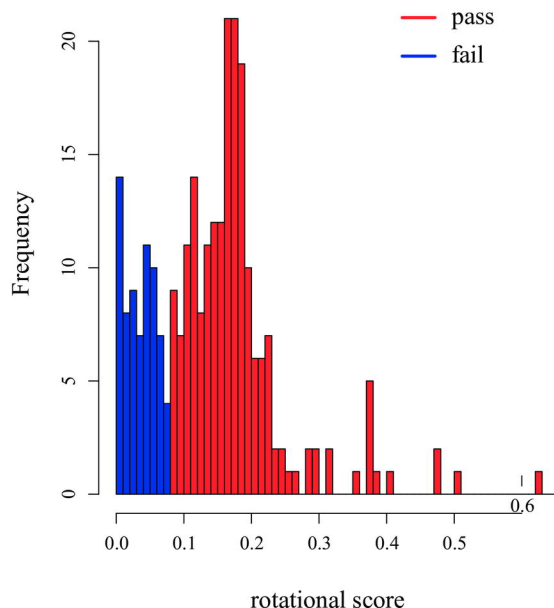
**B**

distribution of nucleosome coverage of historical data



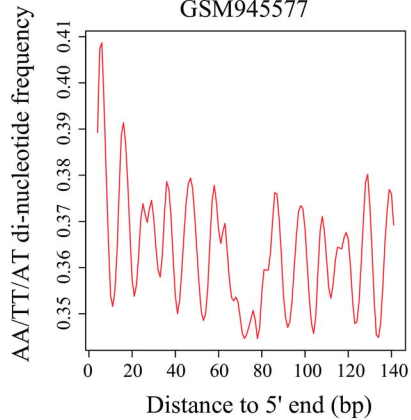
A

distribution of rotational scores of historical data

**B**

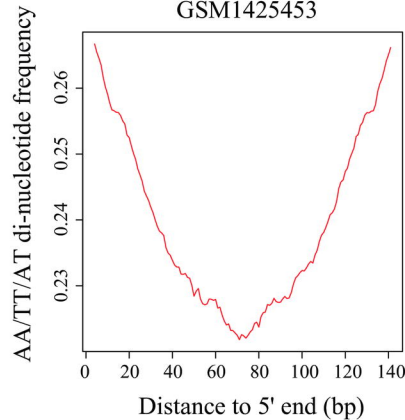
Pass example

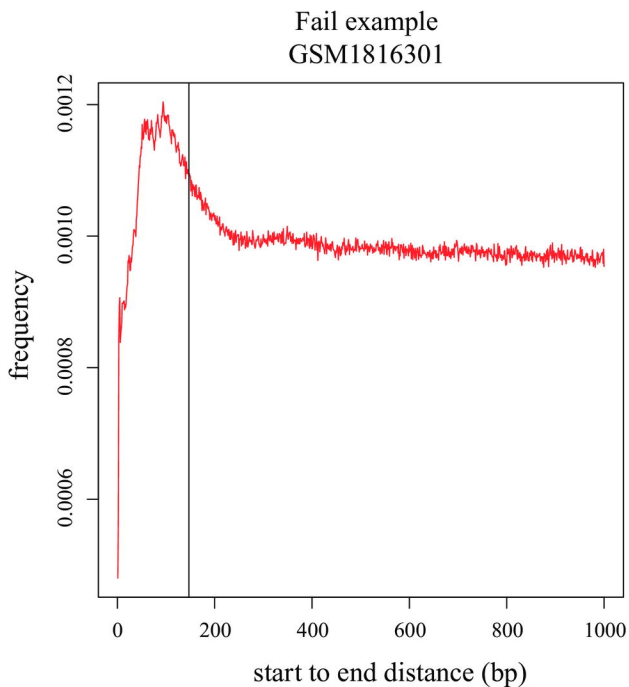
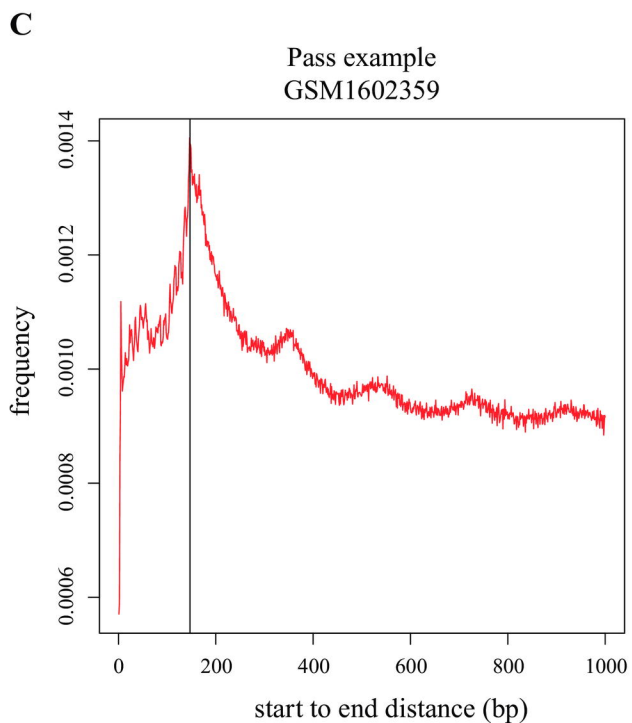
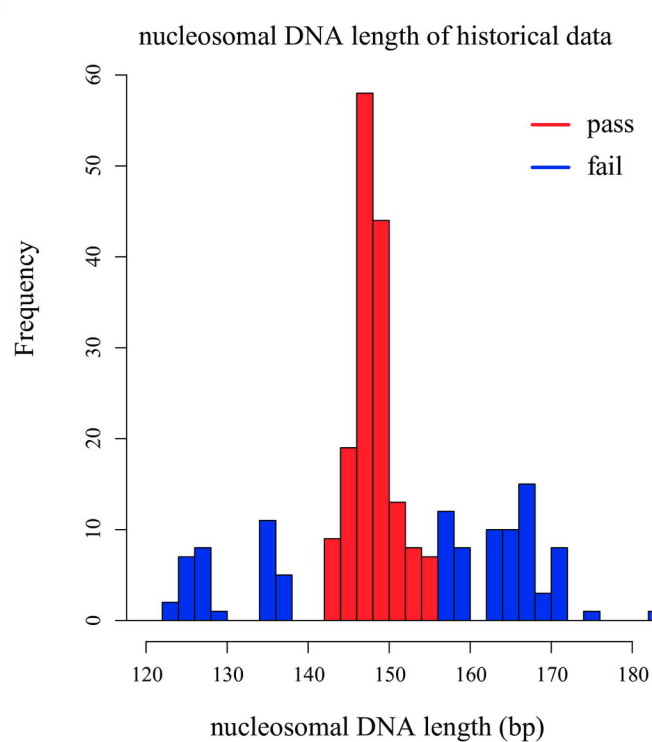
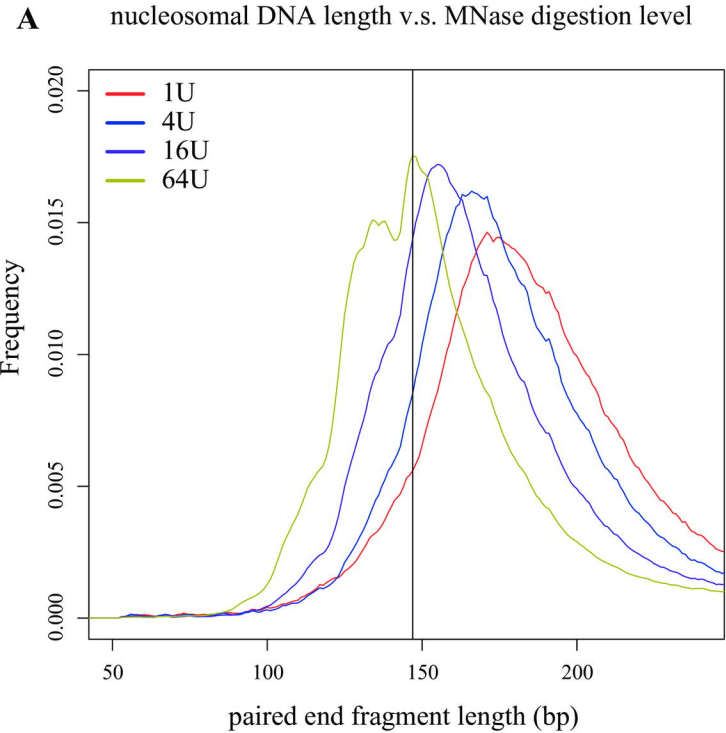
GSM945577



Fail example

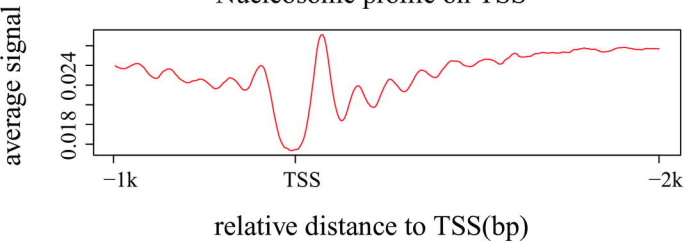
GSM1425453





A

Nucleosome profile on TSS

**B**

Nucleosome profile on custom regions

