

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

**Tn5 transposase-based epigenomic profiling methods
are prone to open chromatin bias**

Meng Wang^{1,2,3} and Yi Zhang^{1,2,3,4,5,*}

¹Howard Hughes Medical Institute, Boston Children’s Hospital, Boston, Massachusetts 02115, USA; ²Program in Cellular and Molecular Medicine, Boston Children’s Hospital, Boston, Massachusetts 02115, USA; ³Division of Hematology/Oncology, Department of Pediatrics, Boston Children’s Hospital, Boston, Massachusetts 02115, USA; ⁴Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA; ⁵Harvard Stem Cell Institute, WAB-149G, 200 Longwood Avenue, Boston, Massachusetts 02115, USA.

*To whom correspondence should be addressed

E-mail: yzhang@genetics.med.harvard.edu

Keywords: Tn5-based methods, low-input epigenomic profiling, open chromatin, comparative analysis

Manuscript information: 16 pages, 3 figures, 3 supplementary figures, 1 supplementary table

28 **Abstract**

29

30 Epigenetic studies of rare biological samples like mammalian oocytes and preimplantation
31 embryos require low input or even single cell epigenomic profiling methods. To reduce sample
32 loss and avoid inefficient immunoprecipitation, several chromatin immuno-cleavage-based
33 methods using Tn5 transposase fused with Protein A/G have been developed to profile histone
34 modifications and transcription factor bindings using small number of cells. The Tn5
35 transposase-based epigenomic profiling methods are featured with simple library construction
36 steps in the same tube, by taking advantage of Tn5 transposase's capability of simultaneous DNA
37 fragmentation and adaptor ligation. However, the Tn5 transposase prefers to cut open chromatin
38 regions. Our comparative analysis shows that Tn5 transposase-based profiling methods are prone
39 to open chromatin bias. The high false positive signals due to biased cleavage in open chromatin
40 could cause misinterpretation of signal distributions and dynamics. Rigorous validation is needed
41 when employing and interpreting results from Tn5 transposase-based epigenomic profiling
42 methods.

43

44

45

46

47

48

49

50

51 **Introduction**

52 Due to the sample loss and inefficient immunoprecipitation of traditional chromatin
53 immunoprecipitation (ChIP)-based methods, low-input epigenomic profiling methods are needed
54 for studying rare samples such as mammalian oocytes and preimplantation embryos¹. Several
55 low-input chromatin immunoprecipitation followed by sequencing (ChIP-seq) methods including
56 ULI-NChIP², scChIP-seq³ and STAR ChIP-seq⁴ have been developed. To overcome inefficient
57 immunoprecipitation, immunoprecipitation-free methods, such as CUT&RUN⁵ and scChIC-seq⁶
58 that use chromatin immuno-cleavage (ChIC)⁷ strategy, have been developed. These low-input
59 methods have been widely used in studying epigenome reprogramming during early embryonic
60 development which have revealed distinct dynamics of different epigenetic markers^{4, 8-12}.

61
62 Recently, the Tn5 transposase-based¹³ library construction is getting popular because it can
63 fragment DNA while simultaneously adding library adaptors thus simplifying experimental
64 procedures and reducing sample loss^{14, 15}. Tn5 transposase-based epigenomic profiling methods
65 utilize Protein A (or Protein G) fused with Tn5 transposase (pA-Tn5) to cleave DNA at the
66 targets of primary antibody, allowing all procedures to be completed in the same tube without
67 immunoprecipitation step, which largely avoided sample loss. Several Tn5-based methods have
68 been developed to capture histone modifications or transcription factor (TF) binding using small
69 number of cells or even single cell, including CUT&Tag¹⁶, CoBATCH¹⁷, ACT-seq¹⁸, itChIP-
70 seq¹⁹, ChIL-seq²⁰ and Stacc-seq²¹. However, the Tn5 transposase is known to prefer accessible
71 DNA regions²². It has been noted that some of the Tn5-based methods are confounded by DNA
72 accessibility²³, but no systematic comparative analysis has been done to determine to what extent
73 the results of these methods are affected by chromatin accessibility. Here we present systematic
74 comparative analyses which reveal that, for some of the methods, overall ~30-50% false positive
75 peaks can be contributed by open chromatin artefacts. Such high level of false positive peaks
76 could affect data interpretation leading to false conclusions, which raises concerns on choosing,
77 developing and interpreting results from Tn5-based epigenomic profiling methods.

78 **Results**

79

80 **Tn5-based epigenomic profiling methods have varied level of biases toward open**

81 **chromatins**

82 CoBATCH¹⁷, CUT&Tag¹⁶, ACT-seq¹⁸ and Stacc-seq²¹ are very similar methods, which are based
83 on the in situ immuno-cleavage strategy (**Fig. 1a**). CoBATCH and CUT&Tag add primary
84 antibodies first, then add the pA-Tn5, while ACT-seq and Stacc-seq pre-incubate primary
85 antibodies with pA-Tn5. However, besides cleavage at the target sites, the free pA-Tn5 has the
86 potential to tagment open chromatins. To wash out free pA-Tn5, these methods employ different
87 washing conditions. The CUT&Tag uses high salt (300 mM NaCl) washing to suppress
88 background tagmentation, CoBATCH and ACT-seq utilize milder washing conditions, while the
89 washing step is optional in Stacc-seq. The itChIP-seq, on the other hand, is an
90 immunoprecipitation-based method that utilizes Tn5 to tagment DNA first before antibodies are
91 added to pull-down the target fragments. Thus, theoretically the itChIP-seq method should not be
92 confounded by open chromatins as antibody-based selectivity is applied after tagmentation.

93

94 To perform systematic comparative analysis of the Tn5-based epigenomic profiling methods, we
95 collected publicly available H3K27me3 data (**Supplementary Table 1**) of mouse embryonic
96 stem cells (mESCs) generated by CoBATCH, CUT&Tag, Stacc-seq and itChIP-seq. Since the
97 protocols for ACT-seq and Stacc-seq are almost identical, and the original ACT-seq study did not
98 include H3K27me3, we analyzed the Stacc-seq data with conclusions applicable to ACT-seq. The
99 bulk ChIP-seq of H3K27me3 in mESCs was used as a reference for comparing the H3K27me3
100 peaks derived from these Tn5-based methods. The two different bulk ChIP-seq datasets^{24, 25} were
101 highly similar (**Supplementary Fig. 1a, b**) and the peaks were considered as true positive peaks
102 in mESCs. To determine whether each Tn5-based method was confounded by open chromatin,
103 we asked whether peaks that were not overlapped with bulk ChIP-seq peaks were instead
104 overlapped with open chromatin peaks derived from ATAC-seq in mESCs. The open chromatin

105 revealed by ATAC-seq²⁶ were highly similar to that revealed by DNase-seq²⁵ in mESCs
106 (**Supplementary Fig. 1c, d**). As a control for low-input epigenomic profiling method without
107 using Tn5 transposase, we included the H3K27me3 dataset in mESCs generated by ULI-NChIP⁸.
108 A genome browser view around the Hoxb locus comparing the signals of Tn5-based methods
109 with those of bulk ChIP-seq, ULI-NChIP, and open chromatin (ATAC-seq and DNase-seq) (**Fig.**
110 **1b**) revealed: 1) CoBATCH, CUT&Tag, and Stacc-seq detected H3K27me3 peaks not present in
111 ChIP-seq or ULI-NChIP but overlapping with ATAC-seq and DNase-seq peaks (shaded); 2) For
112 peaks overlapping with ChIP-seq, the peak patterns were more similar to ATAC-seq and DNase-
113 seq rather than the ChIP-seq; 3) itChIP-seq showed the most similar pattern to that of the ChIP-
114 seq in this region, which was coincident with the fact that the immunoprecipitation-based itChIP-
115 seq procedure is different from the other pA-Tn5 immuno-cleavage-based methods (**Fig. 1a**).

116
117 The above observation raised the possibility that at least some of the Tn5-based methods may be
118 biased toward open chromatin to generate false positive peaks. To explore this possibility, we
119 analyzed the overall signal distribution of these methods by first focusing on the transcription
120 start sites (TSS) of all coding genes. An analysis of the ChIP-seq datasets indicated that the
121 H3K27me3 signals were enriched in the TSSs of a subset of genes consisted of mainly the
122 Polycomb-group (PcG) targets, but with the majority of the genes, mostly of non-PcG targets,
123 lack the H3K27me3 signals around their TSSs (**Fig. 1c**). However, the CoBATCH and Stacc-seq
124 methods detected H3K27me3 enrichment at almost all the TSS regions including the non-PcG
125 targets, which were more similar to the open chromatin patterns detected by ATAC-seq (**Fig. 1c**).
126 The CUT&Tag method detected a weak signal enrichment at the TSSs without H3K27me3 ChIP-
127 seq signals. The itChIP-seq and ULI-NChIP methods detected a pattern more similar to that of
128 ChIP-seq although their signals were generally weaker (**Fig. 1c**). itChIP-seq is not an
129 immunoprecipitation-free method, thus its signals need to be normalized by input control, similar
130 to ChIP-seq and ULI-NChIP. For the immunoprecipitation-free methods, input DNA or IgG
131 control is usually not needed for signal normalization. Nevertheless, we tested whether the

132 confounding of open chromatin signals could be eliminated by using input/IgG control. Using
133 the publicly available input/IgG controls for CoBATCH and Stacc-seq (no input/IgG control for
134 CUT&Tag), we recalculated the H2K27me3 enrichment and found that normalizing with
135 input/IgG did not improve the CoBATCH results (**Fig. 1c**). On the other hand, this normalization
136 did enhance the signals of Stacc-seq overlapping with ChIP-seq, while reduced the signals not
137 overlapping with ChIP-seq (**Fig. 1c**). However, the IgG control normalized Stacc-seq
138 H3K27me3 profile was still more similar to the ATAC-seq profile than that of the H3K27me3
139 ChIP-seq at the non-PcG targets (**Fig. 1c**).

140
141 Next, we focused our analysis on open chromatin regions by dividing the open chromatin regions
142 into two groups that with or without H3K27me3 ChIP-seq signals (**Fig. 1d**). The CoBATCH and
143 Stacc-seq detected signals exhibit a clear enrichment at the open chromatin regions without
144 ChIP-seq signals, and input/IgG control normalization did not change the situation. The
145 CUT&Tag method detected weak signals, while itChIP-seq and ULI-NChIP detected no signals
146 at the open chromatin regions without ChIP-seq signals (**Fig. 1d**). These results indicate that
147 some of the Tn5-based methods, particularly the CoBATCH and Stacc-seq, are biased toward
148 open chromatin peaks that do not have H3K27me3 ChIP-seq signals.

149 **Open chromatin is the source of false positive peaks detected by Tn5-based methods**

150
151 Next we performed quantitative analysis to determine the level that each of the Tn5-based
152 method is confounded by open chromatin. To this end, we used the same criteria (p -value $< 1e-4$
153 and q -value < 0.01) in peak calling for each method. Peaks that overlapped with ChIP-seq peaks
154 were considered as true positives. Peaks that did not overlap with ChIP-seq peaks were
155 considered as potential false positive signals, and were further analyzed to determine whether
156 they could be mapped to open chromatin (**Fig. 2**). We found 5,189 out of the 9,125 CoBATCH
157 peaks did not overlap with ChIP-seq peaks, but were mapped to open chromatin with strong
158 correlation to ATAC-seq signals (**Fig. 2a**). Indeed, 82.3% (4,270 out of 5,189) of the CoBATCH

159 peaks that did not overlap with ChIP-seq peaks were overlapped with ATAC-seq peaks (**Fig. 2b**).
160 Peaks derived from CoBATCH normalized by IgG control showed even more false positives and
161 ATAC-seq signals also enriched in these false positive peaks (**Supplementary Fig. 2a**). A similar
162 analysis of the CUT&Tag dataset revealed 1,387 out of the 6,805 peaks were not overlapped
163 with ChIP-seq peaks (**Fig. 2c**). Of these non-overlapping peaks, 24.2% (335 out of 1,387) were
164 overlapped with ATAC-seq peaks (**Fig. 2d**). For Stacc-seq, 1,687 out of the 6,190 peaks were not
165 overlapped with ChIP-seq peaks, but showed strong open chromatin signals (**Fig. 2e**). Of these
166 non-overlapping peaks, 75.6% (1,275 out of 1,687) were overlapped with ATAC-seq peaks (**Fig.**
167 **2f**). Peaks derived from Stacc-seq normalized by IgG control still showed open chromatin
168 enrichment for the false positive peaks (**Supplementary Fig. 2b**). On the other hand, the itChIP-
169 seq peaks that showed no overlap with ChIP-seq peaks also did not overlap with ATAC-seq
170 peaks (**Fig. 2g, h, Supplementary Fig. 2c**). For ULI-NChIP, although about half of the peak
171 regions showed no ChIP-seq signals, no ATAC-seq signals were detected in these non-
172 overlapping regions (**Fig. 2i, j**). Collectively, these results indicate that the great majority of the
173 non-overlapping peaks detected by CoBATCH or Stacc-seq, and some of the non-overlapping
174 peaks detected by CUT&Tag are mapped to open chromatin regions, and they could be false
175 positive artefacts.

176

177 **High false positive rate due to open chromatin affected global distribution of peaks**

178 To determine the relative reliability of the different Tn5-based epigenomic profiling methods, we
179 next calculated the false positive rate (FPR) caused by the Tn5 bias toward open chromatin. The
180 FPR is calculated by the number of peaks not overlapping with ChIP-seq but overlapping with
181 ATAC-seq peaks, divided with the total number of peaks (**Fig. 3a, Supplementary Fig. 3**). The
182 FPR for CoBATCH was as high as 46.8-54.3%, in contrast the FPR for CUT&Tag was 4.9-5.8%.
183 For Stacc-seq, its FPR was 20.6-35.9%. The itChIP-seq and ULI-NChIP were almost not
184 affected by open chromatin artefacts. Since the FPRs calculated here only considered the non-
185 overlapping peaks, and because the overlapping peaks could also be generated from open

186 chromatin, instead of real H3K27me3 peaks as exemplified in **Fig. 1b**, the FPRs presented here
187 represented the lower limit. To calculate the FPR without a fixed p-value or q-value cutoff for the
188 peaks, we assessed the FPRs for top peaks ranked by p-values for each method. Results shown in
189 **Fig. 3b** indicated that most of the top peaks in CoBATCH represented open chromatin signals.
190 Interestingly, while replicate 1 of the Stacc-seq showed lower FPR for the top peaks but the FPR
191 gradually increased with more peaks, replicate 2 of Stacc-seq showed a relatively consistent FPR
192 (**Fig. 3b**). Consistent with the high FPR, clustering analysis indicated that CoBATCH and Stacc-
193 seq globally resembled open chromatin signals more closely than the H3K27me3 signals (**Fig.**
194 **3c**). These results indicate that CoBATCH and Stacc-seq have high false positive rates and thus
195 great care should be taken in interpreting the data generated by these two methods.

196

197

198 **Discussion**

199

200 In summary, our analysis reveals that the Tn5-based epigenomic profiling methods could capture
201 substantial confounding open chromatin signals. The severity of open chromatin bias varies a lot
202 among the different Tn5-based. CoBATCH and Stacc-seq are prone to open chromatin bias with
203 high false positive rates. Thus, no matter adding the antibody and pA-Tn5 sequentially
204 (CoBATCH) or pre-incubating antibody with pA-Tn5 and adding together (Stacc-seq), both
205 procedures could result in high levels of bias toward open chromatin. Although CUT&Tag
206 showed very weak H3K27me3 signals in non-PcG targets that resembled open chromatin
207 signals, its overall false positive rate due to open chromatin is much lower than that from the
208 CoBATCH despite both share almost identical experimental procedures (**Fig. 1a**). Thus, stringent
209 washing with high salt before tagmentation employed in the CUT&Tag method must have
210 contribute to the reduced open chromatin artefacts²³, but would affect sites with weak binding.
211 Signal normalization with IgG control could not eliminate the confounding signals from open
212 chromatin for CoBATCH and Stacc-seq, while CUT&Tag without IgG control could have much

213 lower biases. This coincides with the practice that no IgG control is needed for in situ immuno-
214 cleavage-based profiling methods. The immunoprecipitation-based itChIP-seq utilizes Tn5
215 transposase to tagment DNA before immunoprecipitating the target DNA, thus its result is not
216 affected by open chromatin. However, it is not an immunoprecipitation-free method, which
217 requires input control, and its signal-to-noise intensities are not comparable to the immuno-
218 cleavage-based methods when using small number of cells.

219
220 Given that Tn5-based methods are prone to open chromatin bias, cautions should be taken when
221 the Tn5-based epigenomic profiling methods are used. We strongly recommend that evaluation
222 of the confounding open chromatin signals and estimation of the FPR are performed under
223 similar experimental conditions before these methods are used. We also suggest that in the future
224 development of Tn5-based epigenomic profiling methods, repressive marks such as H3K27me3
225 or H3K9me3 should be used in evaluating the confounding open chromatin signals, instead of
226 the active H3K4me3 mark used in the original publications of these methods. Since H3K4me3
227 largely colocalizes with open chromatin in mESCs, even the method mainly captures open
228 chromatin signals, the use of H3K4me3 to evaluate would still show high correlation with bulk
229 H3K4me3 ChIP-seq signals. Finally, cautions should be taken when interpreting data generated
230 by Tn5-based epigenomic profiling methods due to the high FPR of open chromatin artefacts.

231

232

233 **References**

234

- 235 1. Carter, B. & Zhao, K. The epigenetic basis of cellular heterogeneity. *Nature reviews.*
236 *Genetics* **22**, 235-250 (2021).
- 237 2. Brind'Amour, J. et al. An ultra-low-input native ChIP-seq protocol for genome-wide
238 profiling of rare cell populations. *Nature communications* **6**, 6033 (2015).
- 239 3. Rotem, A. et al. Single-cell ChIP-seq reveals cell subpopulations defined by chromatin
240 state. *Nature biotechnology* **33**, 1165-1172 (2015).
- 241 4. Zhang, B. et al. Allelic reprogramming of the histone modification H3K4me3 in early
242 mammalian development. *Nature* **537**, 553-557 (2016).

- 243 5. Skene, P.J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution
244 mapping of DNA binding sites. *Elife* **6** (2017).
- 245 6. Ku, W.L. et al. Single-cell chromatin immunocleavage sequencing (scChIC-seq) to profile
246 histone modification. *Nature methods* **16**, 323-325 (2019).
- 247 7. Schmid, M., Durussel, T. & Laemmli, U.K. ChIC and ChEC; genomic mapping of
248 chromatin proteins. *Molecular cell* **16**, 147-157 (2004).
- 249 8. Liu, X. et al. Distinct features of H3K4me3 and H3K27me3 chromatin domains in pre-
250 implantation embryos. *Nature* **537**, 558-562 (2016).
- 251 9. Zheng, H. et al. Resetting Epigenetic Memory by Reprogramming of Histone
252 Modifications in Mammals. *Molecular cell* **63**, 1066-1079 (2016).
- 253 10. Wang, C. et al. Reprogramming of H3K9me3-dependent heterochromatin during
254 mammalian embryo development. *Nat Cell Biol* **20**, 620-631 (2018).
- 255 11. Chen, Z., Djekidel, M.N. & Zhang, Y. Distinct dynamics and functions of H2AK119ub1
256 and H3K27me3 in mouse preimplantation embryos. *Nature genetics* **53**, 551-563 (2021).
- 257 12. Mei, H. et al. H2AK119ub1 guides maternal inheritance and zygotic deposition of
258 H3K27me3 in mouse embryos. *Nature genetics* **53**, 539-550 (2021).
- 259 13. Reznikoff, W.S. Transposon Tn5. *Annu Rev Genet* **42**, 269-286 (2008).
- 260 14. Adey, A. et al. Rapid, low-input, low-bias construction of shotgun fragment libraries by
261 high-density in vitro transposition. *Genome biology* **11**, R119 (2010).
- 262 15. Picelli, S. et al. Tn5 transposase and tagmentation procedures for massively scaled
263 sequencing projects. *Genome research* **24**, 2033-2040 (2014).
- 264 16. Kaya-Okur, H.S. et al. CUT&Tag for efficient epigenomic profiling of small samples and
265 single cells. *Nature communications* **10**, 1930 (2019).
- 266 17. Wang, Q. et al. CoBATCH for High-Throughput Single-Cell Epigenomic Profiling.
267 *Molecular cell* **76**, 206-216 e207 (2019).
- 268 18. Carter, B. et al. Mapping histone modifications in low cell number and single cells using
269 antibody-guided chromatin tagmentation (ACT-seq). *Nature communications* **10**, 3747
270 (2019).
- 271 19. Ai, S. et al. Profiling chromatin states using single-cell itChIP-seq. *Nat Cell Biol* **21**, 1164-
272 1172 (2019).
- 273 20. Harada, A. et al. A chromatin integration labelling method enables epigenomic profiling
274 with lower input. *Nat Cell Biol* **21**, 287-296 (2019).
- 275 21. Liu, B. et al. The landscape of RNA Pol II binding reveals a stepwise transition during
276 ZGA. *Nature* **587**, 139-144 (2020).
- 277 22. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of
278 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-
279 binding proteins and nucleosome position. *Nature methods* **10**, 1213-1218 (2013).
- 280 23. Kaya-Okur, H.S., Janssens, D.H., Henikoff, J.G., Ahmad, K. & Henikoff, S. Efficient low-
281 cost chromatin profiling with CUT&Tag. *Nature protocols* **15**, 3264-3283 (2020).
- 282 24. Perino, M. et al. MTF2 recruits Polycomb Repressive Complex 2 by helical-shape-
283 selective DNA binding. *Nature genetics* **50**, 1002-1010 (2018).

- 284 25. Encode Project Consortium et al. Expanded encyclopaedias of DNA elements in the human
285 and mouse genomes. *Nature* **583**, 699-710 (2020).
- 286 26. Wu, J. et al. The landscape of accessible chromatin in mammalian preimplantation embryos.
287 *Nature* **534**, 652-657 (2016).
- 288

289 **Methods**

290

291 **Data collection.** In the original papers that described each of the Tn5-based method, most used
292 H3K4me3 and/or H3K27me3 in mESCs for validation. However, the H3K4me3 in mESCs is
293 mainly located at promoters in open chromatin regions. It is almost impossible to discriminate
294 the peaks generated by true H3K4me3 or open chromatin. Thus, we used the repressed marker
295 H3K27me3 in mESCs to evaluate the Tn5-based epigenomic profiling methods (summarized in
296 Supplementary Table 1), which had the most publicly available datasets for different methods
297 besides H3K4me3. For multiple replicates with the same or different number of cells for each
298 method, we used the one with the best signal-to-noise ratio as the representative result for each
299 method.

300

301 **Peak calling and signal track generation.** For ChIP-seq, ULI-NChIP and Tn5-based methods,
302 raw sequencing reads were first trimmed using Trimmomatic²⁷ (version 0.39) to remove
303 sequencing adaptors and low-quality reads. The cleaned reads were mapped to mm10 reference
304 genome using bowtie2²⁸ (version 2.4.2) with parameters: --local --very-sensitive-local --no-unal -
305 -no-mixed --no-discordant --dovetail -I 10 -X 700 --soft-clipped-unmapped-tlen. PCR duplicates
306 were removed with Picard MarkDuplicates (version 2.23.4). Reads with mapping quality at least
307 30 were kept. For Tn5-based methods, proper paired reads with fragment length at least 178bp
308 (nucleosome DNA size 140bp + 2 × Tn5 steric hindrance 19bp at both sides) were kept. For Tn5-
309 based methods, to increase peaks resolution, the start and end positions for each fragment (one
310 read pair) were shifted for 19bp toward internal to account for the steric hindrance of Tn5
311 enzyme. Peaks were called using MACS2²⁹ callpeak (version 2.2.7.1) with parameters: -B -
312 SPMR -p 1e-4 -g mm --broad --broad-cutoff 1e-4 --keep-dup all --scale-to large. Peaks were
313 further filtered with q-value<0.01. The fold-change signal tracks were generated using MACS2
314 bdgcmp with input of treat-pileup and control-lambda bedgraph files generated from MACS2
315 callpeak in the last step. Peaks overlapping with mm10 blacklist regions

316 (<https://www.encodeproject.org/files/ENCFF547MET/>) were removed. The ChIP-seq results
317 were pooling of two replicates. The ULI-NChIP results were pooling of all four replicates. The
318 ATAC-seq and DNase-seq datasets were analyzed using ENCODE ATAC-seq pipeline (version
319 1.9.3, <https://github.com/ENCODE-DCC/atac-seq-pipeline>).

320
321 **Peak comparison.** Peaks were compared using bedtools³⁰ intersect (version 2.29.2). Peaks with
322 at least half-length intersecting with ChIP-seq / ATAC-seq / DNase-seq peaks were considered as
323 overlapping (bedtools intersect parameters for getting overlapping peaks: -u -f 0.5; parameters
324 for getting non-overlapping peaks: -v -f 0.5). The signal enrichment heatmaps were plotted using
325 deeptools³¹ (version 3.5.0) computeMatrix and plotHeatmap. The TSSs of coding genes in the
326 mouse genome were from GENCODE³² mouse gene set M24. The genome browser snapshot
327 was generated with R package karyoploteR³³ (version 1.18.0). The Pearson correlation and
328 clustering analysis were performed using deeptools multiBigwigSummary and plotCorrelation
329 with 5kb bin size and outliers removed.

330

331 **Data availability**

332 The public datasets used in this study are summarized in Supplementary Table 1.

333

334 **Code availability**

335 The code used to analyze the sequencing data is available at GitHub:

336 <https://github.com/YiZhang-lab/ChIPpipes>

337

338

339 **References**

- 340 27. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
341 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 342 28. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nature methods*
343 **9**, 357-359 (2012).
- 344 29. Zhang, Y. et al. Model-based analysis of ChIP-Seq (MACS). *Genome biology* **9**, R137

- 345 (2008).
346 30. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic
347 features. *Bioinformatics* **26**, 841-842 (2010).
348 31. Ramirez, F. et al. deepTools2: a next generation web server for deep-sequencing data
349 analysis. *Nucleic acids research* **44**, W160-165 (2016).
350 32. Frankish, A. et al. GENCODE reference annotation for the human and mouse genomes.
351 *Nucleic acids research* **47**, D766-D773 (2019).
352 33. Gel, B. & Serra, E. karyoploteR: an R/Bioconductor package to plot customizable genomes
353 displaying arbitrary data. *Bioinformatics* **33**, 3088-3090 (2017).

354

355

356

357 **Acknowledgements**

358 We thank Drs. Chunxia Zhang and Zhiyuan Chen for discussion and critical reading of the
359 manuscript. This project was supported by the NIH (R01HD092465) and the HHMI. Y.Z. is an
360 Investigator of the Howard Hughes Medical Institute.

361

362 **Author contributions**

363 Y.Z. supervised the project. M.W. performed the analysis. M.W. and Y.Z. wrote the manuscript.

364

365 **Competing interests**

366 The authors declare no competing interests.

367

368

369

370

371 **Figure legends**

372

373 **Fig. 1 | Signal distributions of Tn5-based epigenomic profiling methods**

374 **a**, Major experimental procedures for different Tn5-based epigenomic profiling methods. pA-
375 Tn5: Protein A and Tn5 fusion complex; Ab: primary antibody.

376 **b**, Genome browser snapshot around Hoxb cluster in mESCs for open chromatin fold-change
377 signals (ATAC-seq and DNase-seq), and H3K27me3 fold-change signals for two ChIP-seq
378 datasets, Tn5-based methods and ULI-NChIP. The itChIP-seq fold-change signals were
379 normalized by input. The signals for CoBATCH, CUT&Tag and Stacc-seq were fold-changes
380 over background.

381 **c**, H3K27me3 signal enrichment for different methods around the transcription start sites
382 (TSS±2kb) of mouse coding genes, and was compared to ATAC-seq signals around the TSSs
383 (FC: fold-change over background/input).

384 **d**, Signal enrichment for different methods at all open chromatin regions (n1: open chromatin
385 regions with H3K27me3 ChIP-seq signals; n2: open chromatin regions without H3K27me3
386 ChIP-seq signals; C: center of ATAC-seq peaks; FC: fold-change over background/input).

387

388 **Fig. 2 | Evaluation of peaks from Tn5-based epigenomic profiling methods**

389 Significant peaks (p-value<1e-4 and q-value<0.01) called from each method (**a**: CoBATCH, **c**:
390 CUT&Tag, **e**: Stacc-seq, **g**: itChIP-seq, **i**: ULI-NChIP) were divided into two parts: n1 – peaks
391 overlapping with ChIP-seq peaks, n2 – peaks not overlapping with ChIP-seq peaks, and
392 compared with open chromatin signals measured by ATAC-seq (C: center of peaks called from
393 each method; FC: fold-change over background/input). The itChIP-seq results shown in **g** and **h**
394 were from 10k cells. For each method (**b**: CoBATCH, **d**: CUT&Tag, **f**: Stacc-seq, **h**: itChIP-seq,
395 **j**: ULI-NChIP), peaks that were not overlapped with ChIP-seq peaks were further compared to
396 ATAC-seq peaks (+ATAC: overlapping with ATAC-seq peaks; -ATAC: not overlapping with
397 ATAC-seq peaks).

398 **Fig. 3 | False positive rates of Tn5-based methods due to open chromatin**

399 **a**, Overall false positive rate (FPR) due to open chromatin (measured by ATAC-seq) artefacts for
400 each method. The number of cells (#cell) used for each library was indicated below each bar.

401 **b**, False positive rate due to open chromatin (measured by ATAC-seq) artefacts for the top peaks
402 in each method.

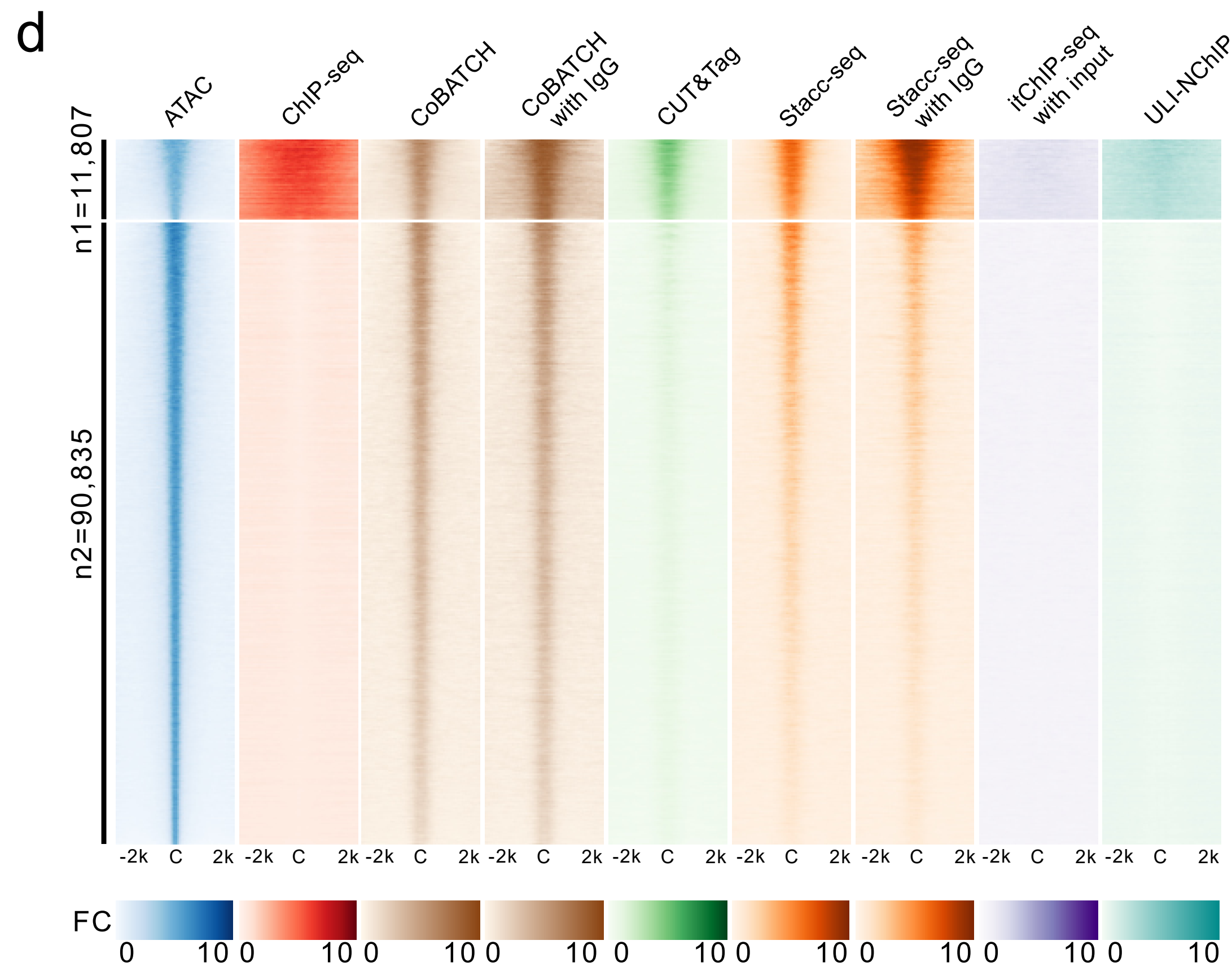
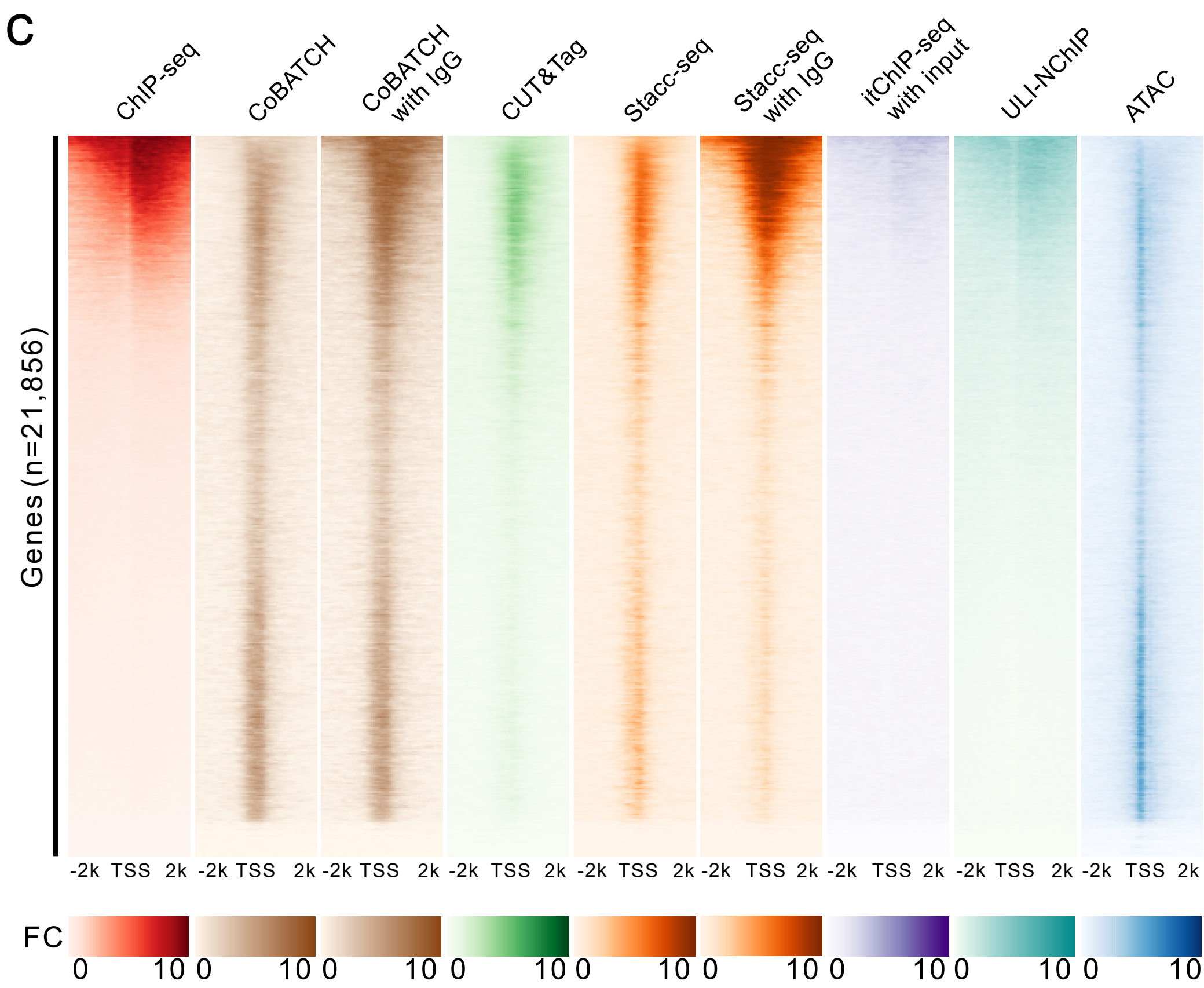
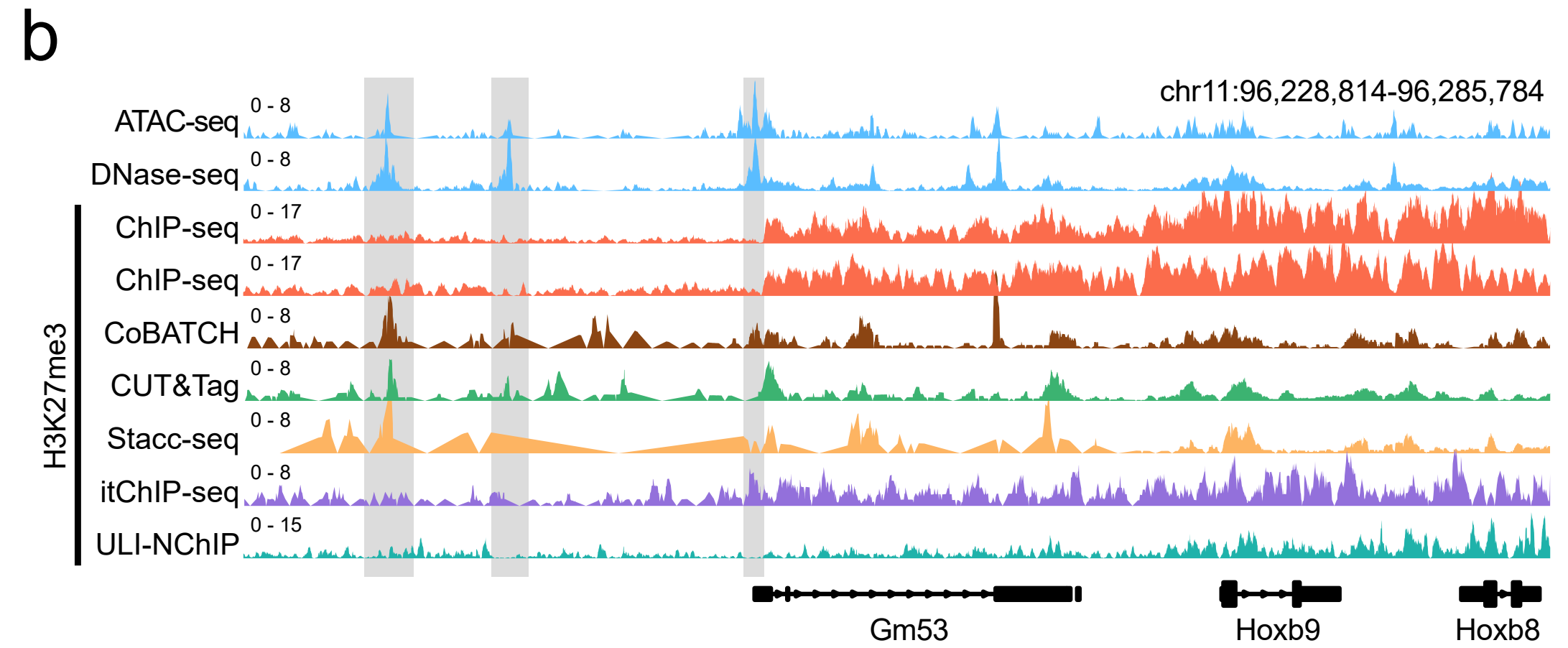
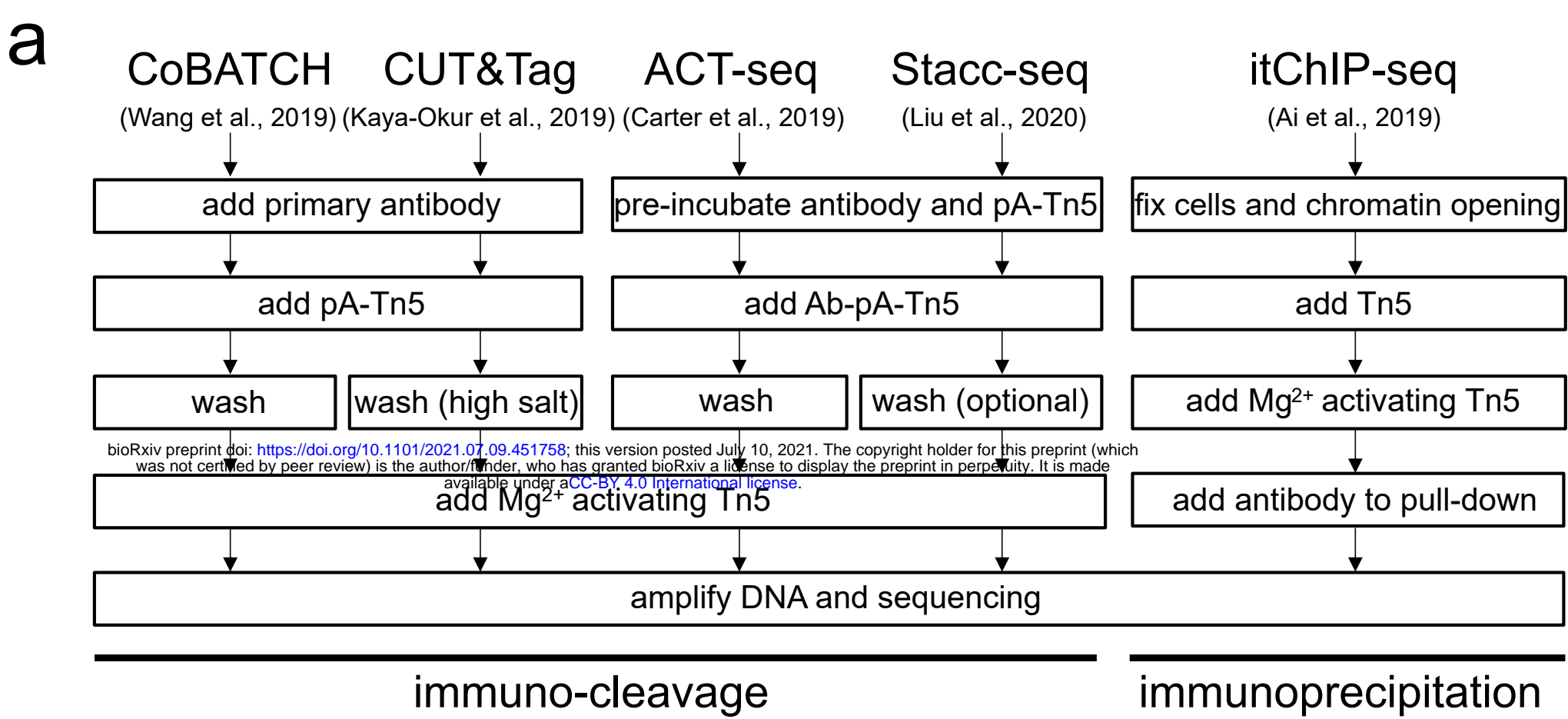
403 **c**, Clustering of global H3K27me3 signals of each method with ATAC-seq and H3K27me3 ChIP-
404 seq based on the Pearson correlation between any two methods (The Pearson correlation
405 coefficients were shown in each box; bin size: 5kb).

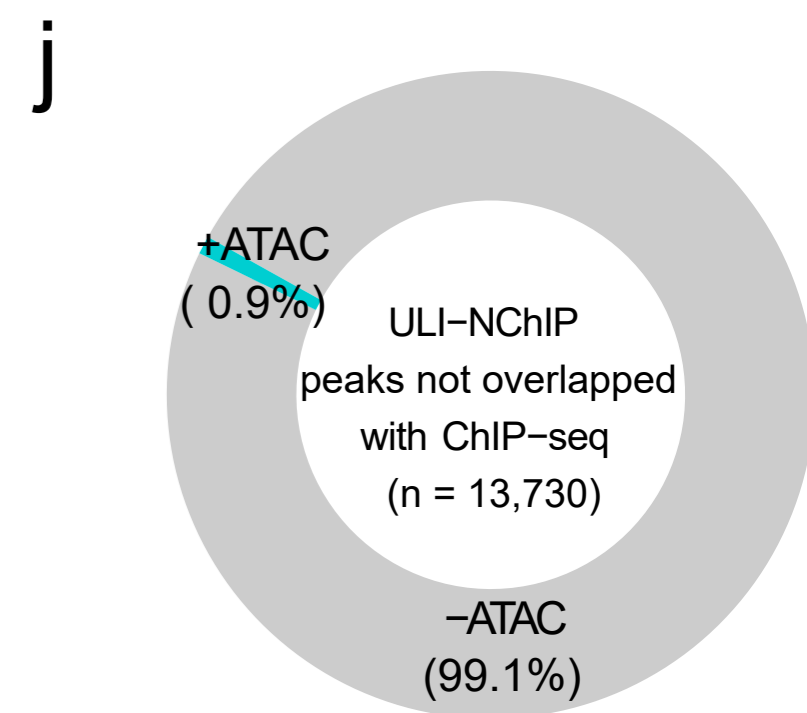
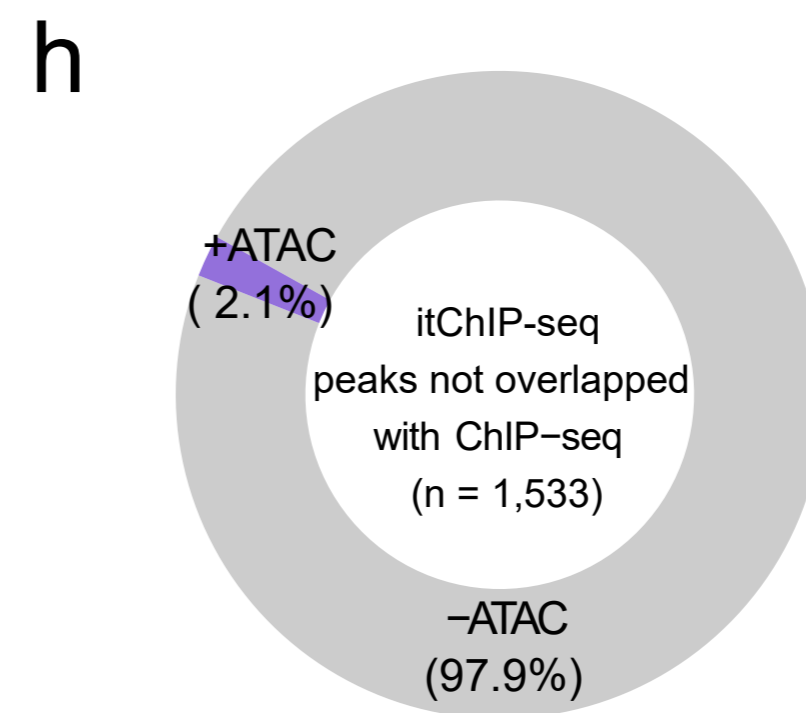
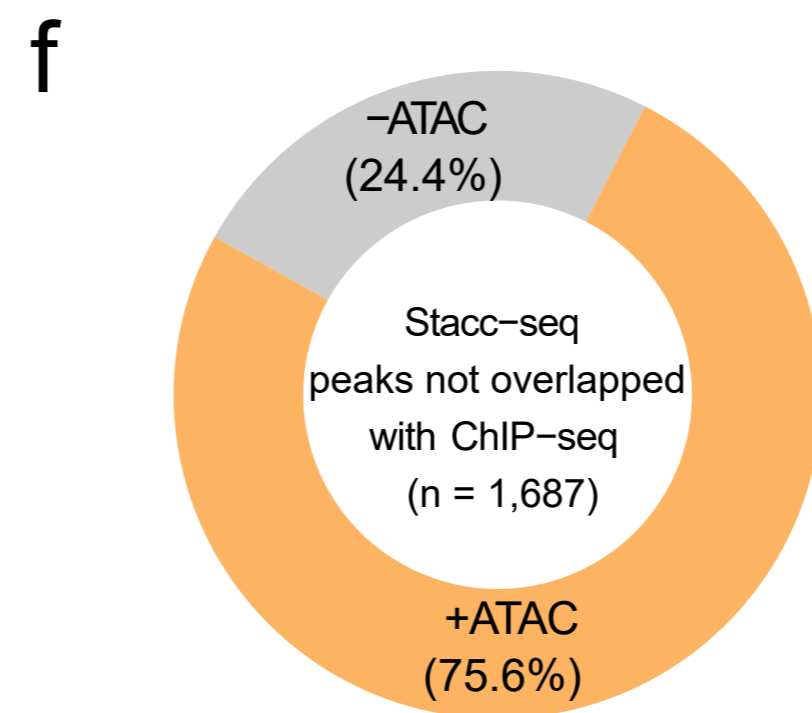
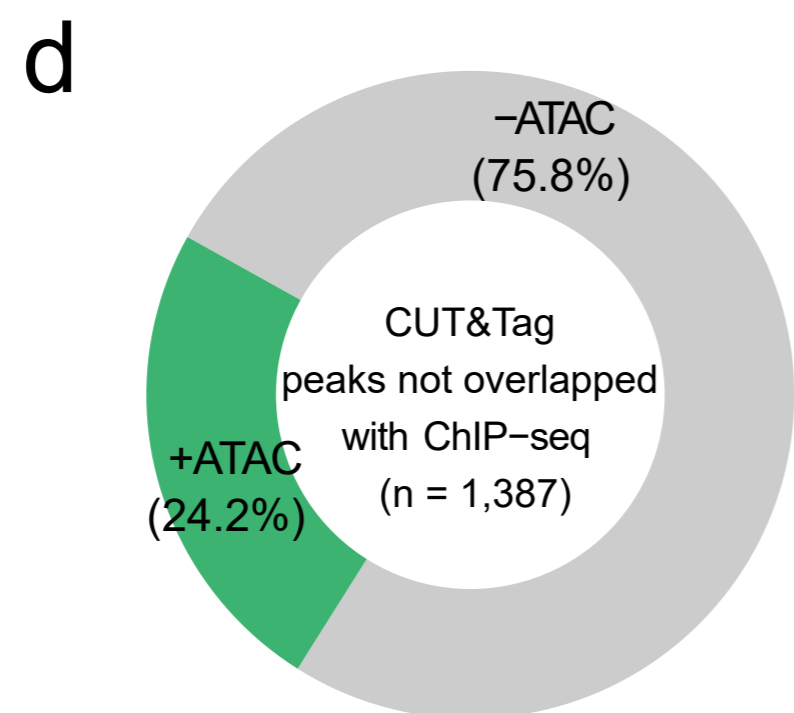
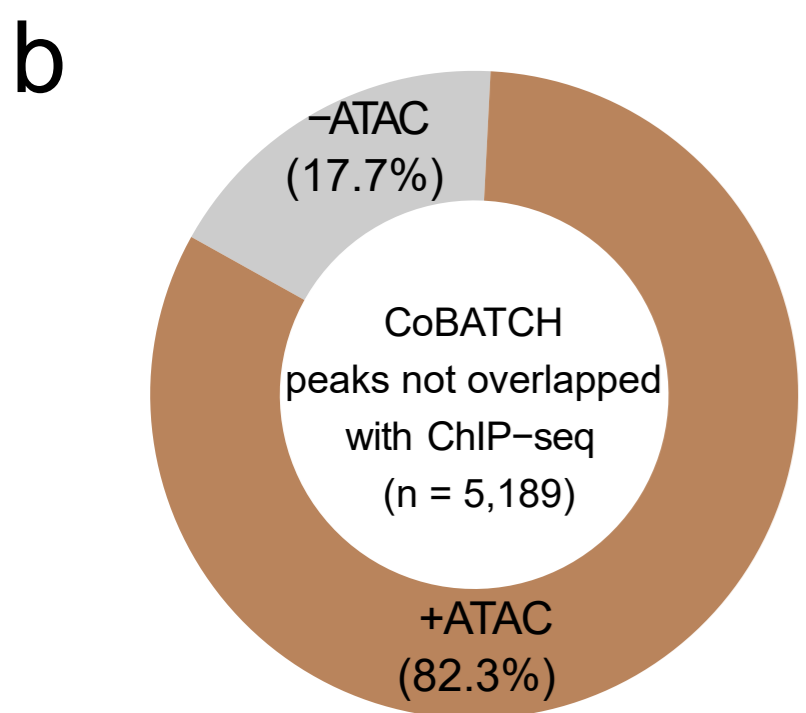
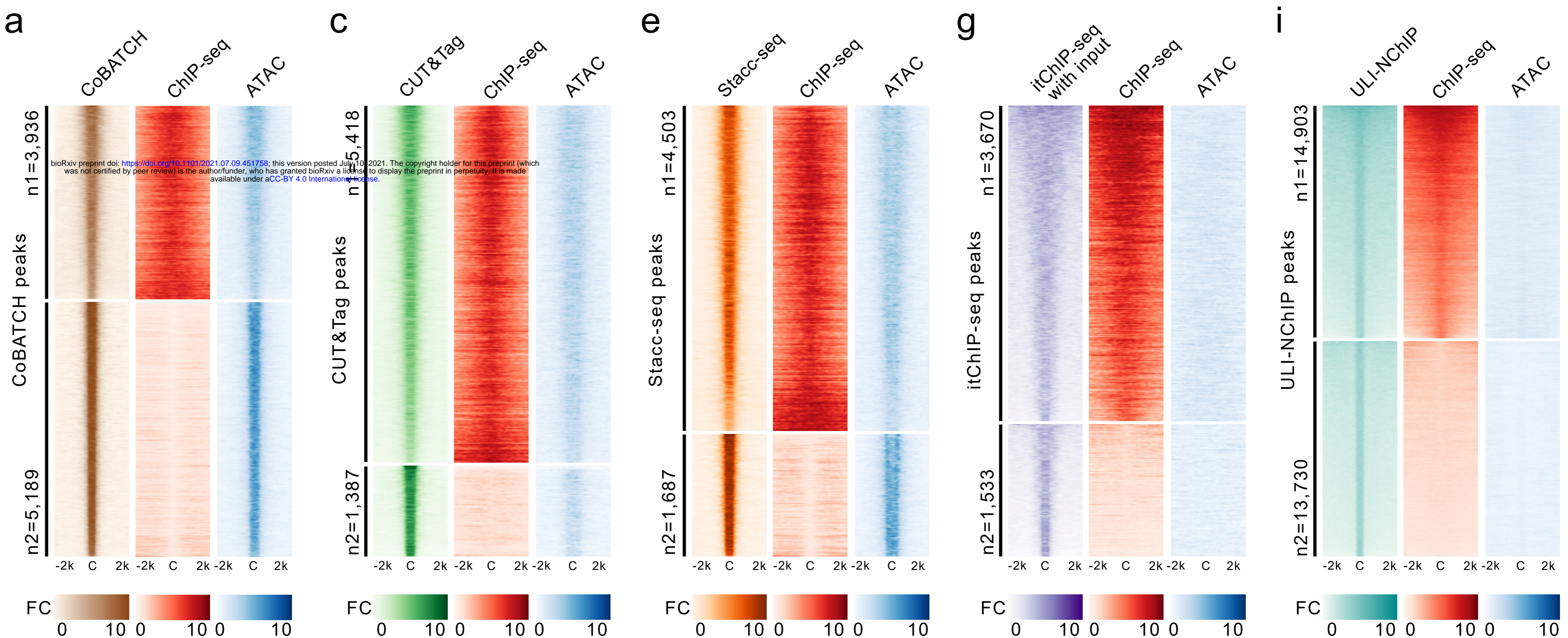
406

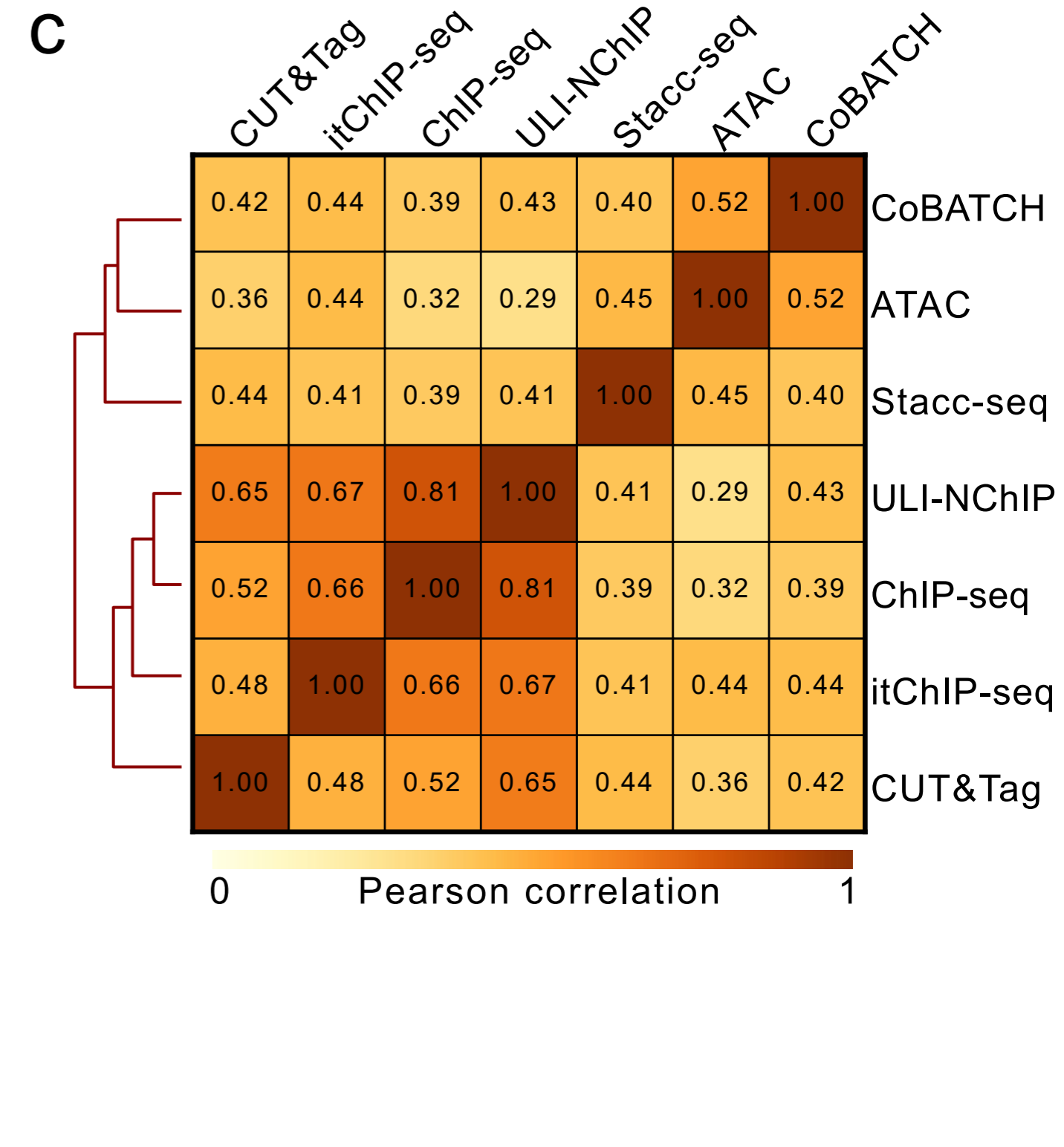
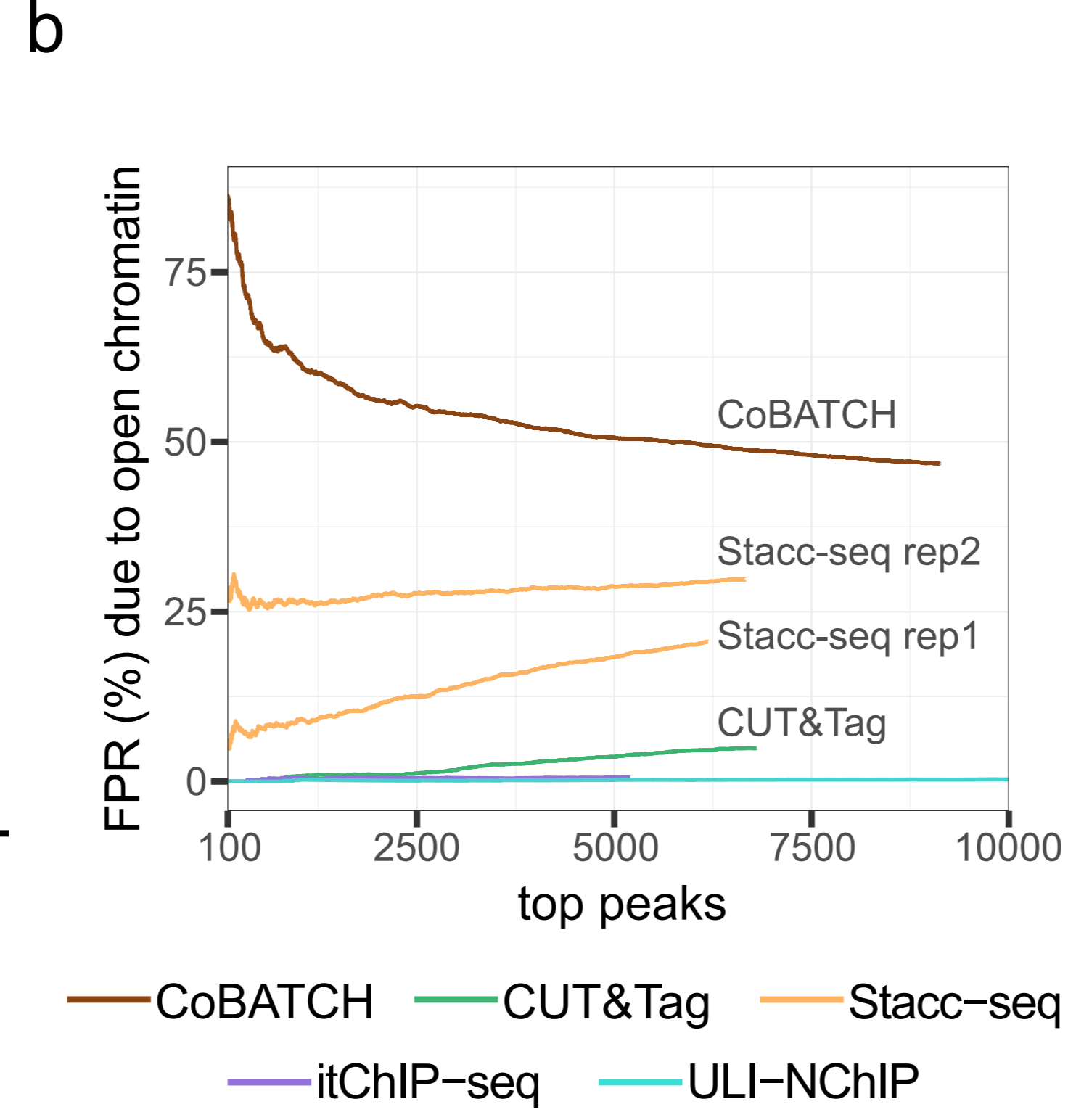
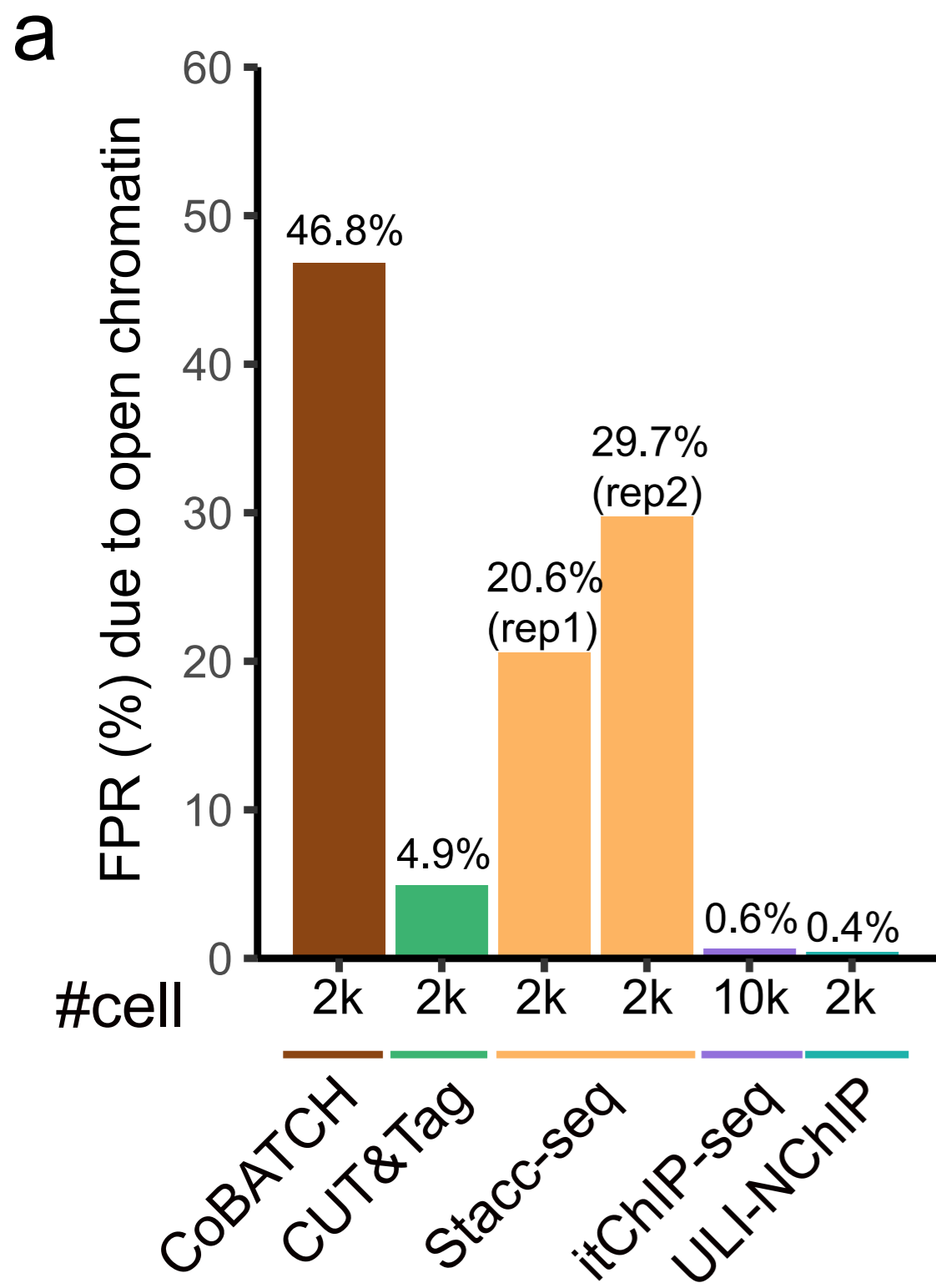
407

408

409







Supplementary for

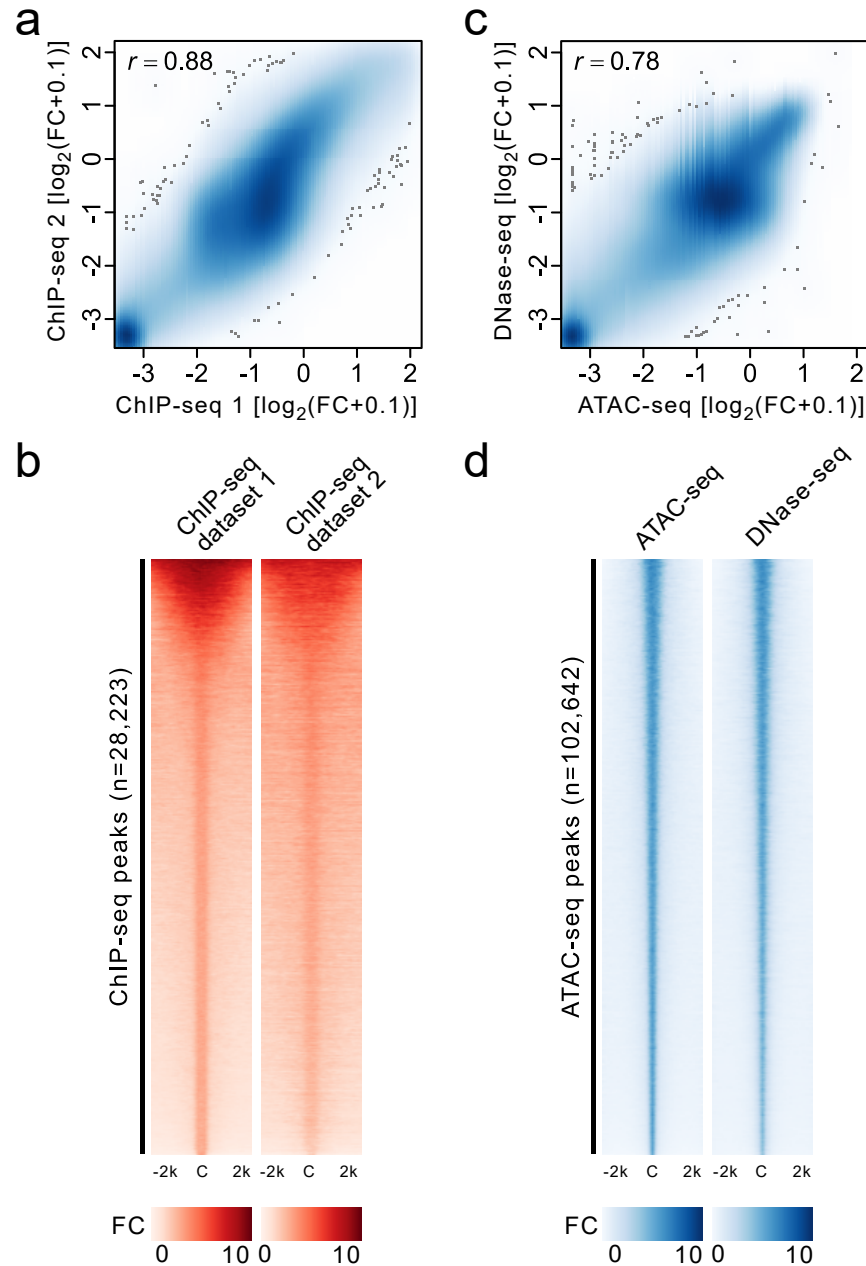
Tn5 transposase-based epigenomic profiling methods are prone to open chromatin bias

Meng Wang^{1,2,3} and Yi Zhang^{1,2,3,4,5,*}

¹Howard Hughes Medical Institute, Boston Children's Hospital, Boston, Massachusetts 02115, USA; ²Program in Cellular and Molecular Medicine, Boston Children's Hospital, Boston, Massachusetts 02115, USA; ³Division of Hematology/Oncology, Department of Pediatrics, Boston Children's Hospital, Boston, Massachusetts 02115, USA; ⁴Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA; ⁵Harvard Stem Cell Institute, WAB-149G, 200 Longwood Avenue, Boston, Massachusetts 02115, USA.

*To whom correspondence should be addressed

E-mail: yizhang@genetics.med.harvard.edu



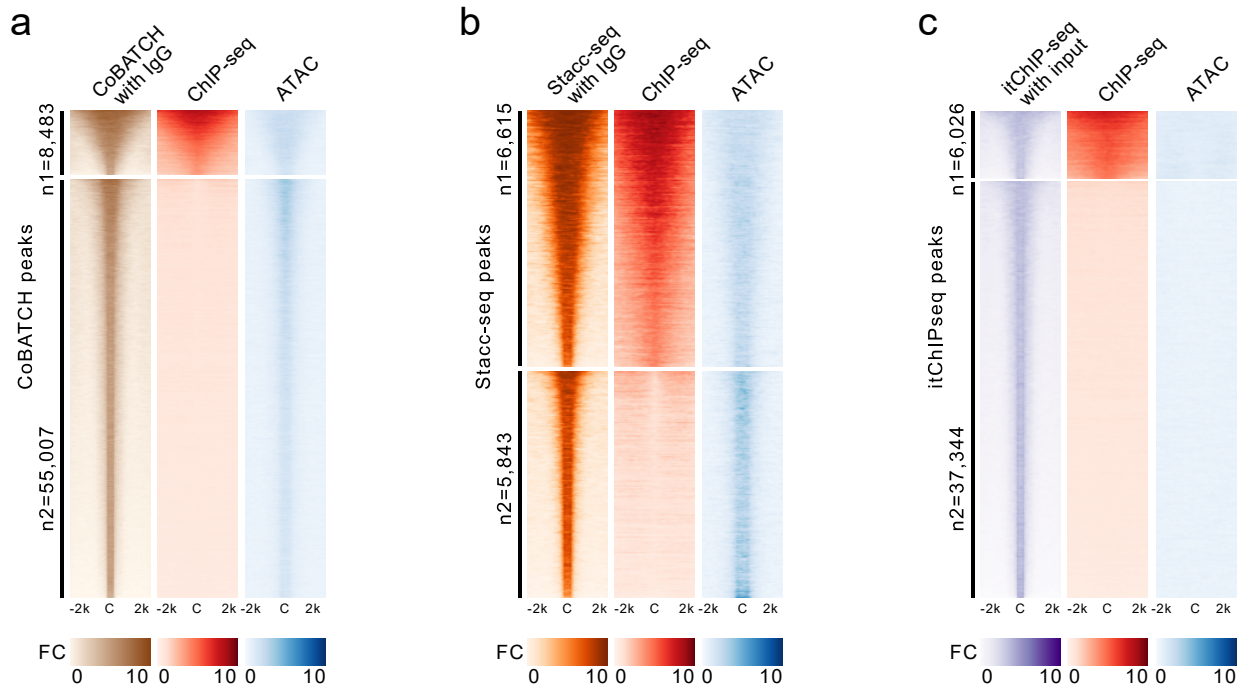
Supplementary Fig. 1 | Different ChIP-seq and ATAC-seq datasets in mESC are consistent

a, Pearson correlation of two H3K27me3 ChIP-seq datasets in mESC (bin size: 5kb).

b, Heatmap comparing H3K27me3 fold-change (FC) signals at the ChIP-seq peaks of two ChIP-seq datasets in mESCs (C: center of peaks in ChIP-seq dataset 1).

c, Pearson correlation of ATAC-seq and DNase-seq in mESC (bin size: 5kb).

d, Heatmap comparing open chromatin fold-change (FC) signals measured by ATAC-seq and DNase-seq in mESCs (C: center of peaks in ATAC-seq).

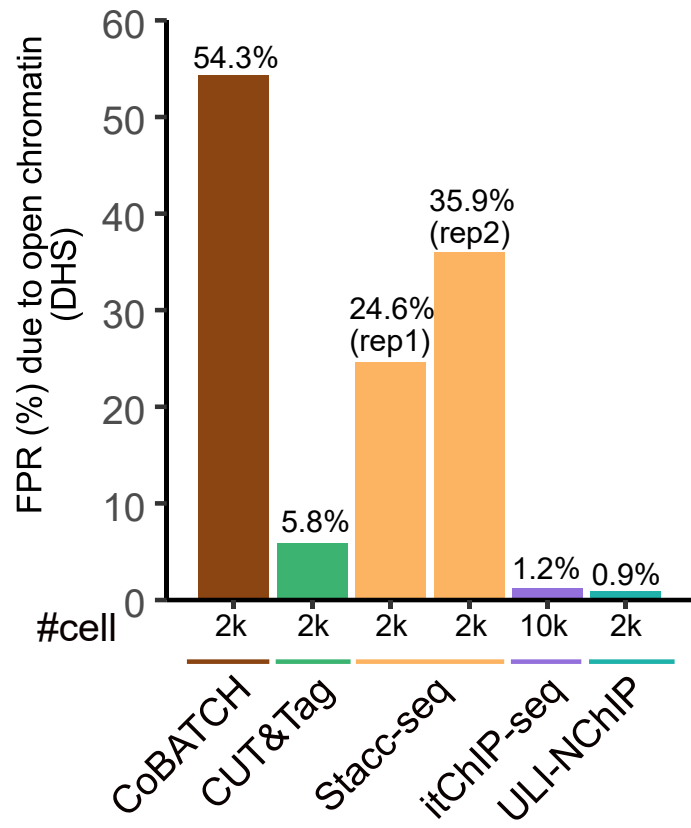


Supplementary Fig. 2 | Evaluation of peaks called from different methods with input / IgG normalization

a, Significant peaks (q-value<0.01) called from CoBATCH with IgG control normalization were compared with open chromatin signals measured by ATAC-seq (n1: CoBATCH peaks overlapping with ChIP-seq peaks; n2: CoBATCH peaks not overlapping with ChIP-seq peaks; C: center of CoBATCH peaks; FC: fold-change over IgG control).

b, Significant peaks (q-value<0.01) called from Stacc-seq with IgG control normalization were compared with open chromatin signals measured by ATAC-seq (n1: Stacc-seq peaks overlapping with ChIP-seq peaks; n2: Stacc-seq peaks not overlapping with ChIP-seq peaks; C: center of Stacc-seq peaks; FC: fold-change over IgG control).

c, Significant peaks (q-value<0.01) called from itChIP-seq (100 cells) with input control normalization were compared with open chromatin signals measured by ATAC-seq (n1: itChIP-seq peaks overlapping with ChIP-seq peaks; n2: itChIP-seq peaks not overlapping with ChIP-seq peaks; C: center of itChIP-seq peaks; FC: fold-change over input control).



Supplementary Fig. 3 | Overall false positive rate (FPR) due to open chromatin (measured by DNase-seq) artefacts for each method. The number of cells (#cell) used for each library was indicated below each bar.

Supplementary Table 1 | Summary of public datasets used in this study.

Method	Sample name	Cell number	Data source	Accession
CoBATCH	2k_mESC	2,000	GEO	GSM3711220
	100_mESC_rep1	100	GEO	GSM3711218
	100_mESC_rep2	100	GEO	GSM3711219
	IgG_2k_mESC_rep1	2,000	GEO	GSM3893775
	IgG_2k_mESC_rep2	2,000	GEO	GSM3893776
CUT&Tag	2k_mESC_rep1	2,000	GEO	GSM4476407
	2k_mESC_rep2	2,000	GEO	GSM4476406
Stacc-seq	2k_mESC_rep1	2,000	GEO	GSM4010607
	2k_mESC_rep2	2,000	GEO	GSM4010608
	IgG_mESC	NA	GEO	GSM4010609
itChIP-seq	10k_mESC_rep1	10,000	GEO	GSM3609659
	10k_mESC_rep2	10,000	GEO	GSM3609660
	100_mESC_rep1	100	GEO	GSM3609661
	100_mESC_rep2	100	GEO	GSM3609662
	input	10,000	GEO	GSM3609658
ULI-NChIP	500_mESC_rep1	500	GEO	GSM2082708
	500_mESC_rep2	500	GEO	GSM2082709
	500_mESC_rep3	500	GEO	GSM2082710
	500_mESC_rep4	500	GEO	GSM2082711
	input	500	GEO	GSM2082705
ChIP-seq (dataset 1)	mESC_rep1	bulk	GEO	GSM2472743
	mESC_rep2	bulk	GEO	GSM2472744
	input	bulk	GEO	GSM2472755
ChIP-seq (dataset 2)	mESC_rep1	bulk	ENCODE	ENCFF001ZIB
	mESC_rep2	bulk	ENCODE	ENCFF001ZIH
	input_rep1	bulk	ENCODE	ENCFF001ZGK
	input_rep2	bulk	ENCODE	ENCFF001ZGM
ATAC-seq	50k_mESC	50,000	GEO	GSM2156965
DNase-seq	bulk	bulk	ENCODE	ENCSR000CMW