1  # Sensitive and robust assessment of ChIP-seq read

2  # distribution using a strand-shift profile

3  Ryuichiro Nakato[1,*] and Katsuhiko Shirahige[1]

4  [1]Institute of Molecular and Cellular Biosciences, The University of Tokyo, 1-1-1 Yayoi,

5  Bunkyo-ku, Tokyo 113-0032, JAPAN

6  *To whom correspondence should be addressed: rnakato@iam.u-tokyo.ac.jp

7

8  **ABSTRACT**

9  Chromatin immunoprecipitation followed by sequencing (ChIP-seq) can detect read-enriched

10  DNA loci for point-source (e.g., transcription factor binding) and broad-source factors (e.g.,

11  several histone modifications). Although numerous quality metrics for ChIP-seq data have

12  been developed, the 'peaks' thus obtained are still difficult to assess with respect to signal-

13  to-noise ratio (S/N) especially for broad-source factors, and peak reliability. Here we

14  introduce SSP (strand-shift profile), a tool to assess the quality of ChIP-seq data without

15  peak calling. SSP provides metrics to quantify the S/N for both point- and broad-source

16  factors, and to estimate peak reliability based on the mapped-read distribution throughout a

17  genome. We carried out an in-depth validation of our method using over 1,000 publicly

18  available ChIP-seq datasets, along with virtual data, to demonstrate that SSP is more

19  sensitive than existing tools for both point- and broad-source factors because of the larger

20  dynamic range of the S/N score, and robust for various cell types and sequencing depth. We

21  also found that SSP can identify low-quality samples that cannot be identified by quality

22  metrics currently available. Finally, SSP provides an additional metric to avoid "hidden-

23  duplicate reads" that cause aberrantly high S/Ns in the strand-shift profile. This metric can

24  also contribute to estimation of peak mode (point- or broad-source) of each sample. Our

25    approach provides a useful way to obtain information about sample quality and traits for

26    ChIP-seq analyses.

27

28    **Introduction**

29    Chromatin immunoprecipitation followed by sequencing (ChIP-seq) analysis identifies DNA

30    loci of transcriptional factors (TFs) binding (i.e., point-source) as well as broadly distributed

31    histone modifications (i.e., broad-source) [1, 2]. In a ChIP experiment, immunoprecipitated

32    DNA fragments are sequenced to reads, which are mapped to a reference genome, and

33    statistically significant read enrichments (as compared with a corresponding input sample)

34    are detected as peaks. Large consortia such as ENCODE [3], NIH ROADMAP [4] and IHEC

35    [5] enable us to utilize thousands of ChIP-seq data for diverse cell lines and tissues. To

36    handle such large-scale data, objective quality metrics for quantitative assessment are

37    essential to automatically find samples which should be rejected or require a specific

38    consideration to be included in the analysis. Numerous computational measures for ChIP-

39    seq analysis have been developed, which include read quality, library complexity, and GC

40    content [6, 7]. Despite great effort, however, the current approach for assessing peaks is

41    insufficient.

42    To assess the success of the immunoprecipitation step, signal-to-noise ratio (S/N) is

43    assessed, and the value should be high for ChIP samples and low for input samples. A

44    straightforward way to evaluate the S/N is to count the number of obtained peaks and/or

45    calculate the fraction of reads falling within peak regions (called FRiP), but these ways

46    depend on sequencing depth and peak-calling parameters. In contrast, cross-correlation

47    analysis [6] evaluates the S/N without the need for a peak-calling procedure. It estimates the

48    Pearson correlation coefficient between the read densities mapped on the forward and

49    reverse DNA strands upon shifting from one strand to the other (see Supplemental Fig. S1

50    for an example). Such a "strand-shift profile" typically peaks at the shift corresponding to the

51    DNA fragment length, which increases as the S/N of the sample increases. This tendency

52  has also been used to estimate fragment length from single-end reads. There is also a spike

53  at the read-length shift that arises from repetitive sequences [8]. Based on this observation,

54  cross-correlation analysis calculates two metrics, namely the normalized strand coefficient

55  (NSC) and the relative strand correlation (RSC), which quantify the fragment length peak

56  relative to background level and relative to the read length peak, respectively (see Results,

57  "Method overview", for details). These metrics have been used in the ENCODE, ROADMAP

58  and IHEC consortia. A strand-shift profile strategy based on the Hamming distance was also

59  proposed for rapid computation (Hansen et al. 2015). Whereas these tools are useful for

60  point-source factors, broad-source factors (e.g., H3K9me3) often have marginal or truly low

61  scores compared with input samples, even when the samples are of high quality [6].

62  Moreover, these S/N indicators do not evaluate the reliability of obtained peaks, that is,

63  amount of false positives which are derived from read distribution bias (e.g., GC bias) [9].

64  Visual inspection at a limited number of sites is effective but not sufficient to explain the

65  properties of read distribution in a whole genome. Consequently, genome-wide assessment

66  of ChIP-seq peak quality still presents challenges that current protocols cannot circumvent.

67  In this work, we present a new method, SSP, which is based on a strand-shift profile using

68  the Jaccard index to assess S/N, peak reliability and properties of read enrichment in ChIP-

69  seq data. We evaluated the performance of SSP using an extensive dataset of ChIP-seq

70  samples for various cell types obtained from the ENCODE, ROADMAP, and other projects,

71  along with simulated experiments. We demonstrate that SSP provides a more sensitive S/N

72  indicator than current methods both for point- and broad-source marks and is robust for

73  various cell types and sequencing depth. We also found that "hidden-duplicate reads" in a

74  sample confound the strand-shift profile because they cause unexpected enrichment,

75  resulting in calculation of aberrantly high S/Ns. Therefore we additionally developed metrics

76  to overcome this problem, which can also be used to estimate peak mode (point or broad

77  source) of each sample. SSP provides a useful way to assess and obtain additional

78  information about sample quality and traits for ChIP-seq analyses.

79

80  **Results**

81  **Method overview**

82  Fig. 1 presents an overview of SSP (see Methods for details). Using mapped reads as input,

83  SSP generates the strand-specific vectors for forward and reverse strands (step 1). Because

84  sequenced reads that are mapped to the same genomic position are removed as duplicate

85  reads [6], each element of a strand-specific vector is binary, that is, either zero (unmapped)

86  or one (mapped). This binary vector can be handled by computationally fast bit operations in

87  C++ [10]. SSP calculates the Jaccard index between binary vectors of forward and reverse

88  strands for each strand shift $d$, which is then normalized by total read number and

89  chromosome length (step 2). The magnitude of the Jaccard score reflects the co-occurrence

90  of reads mapped on the forward and reverse strands with distance $d$. Whereas the Pearson

91  correlation and Hamming distance confer equal weight to pairs of mapped bases (1,1) and

92  unmapped bases (0,0), the Jaccard index focuses on the mapped bases because

93  unmapped bases can often coincide owing to the lack of sequencing depth and low-

94  mappable regions.

95  A strand-shift profile is generated within −500 bp < $d$ < 1 Mbp (step 3). NSC and RSC are

96  then calculated in the same manner as a cross-correlation analysis. Whereas existing

97  methods use ~1,000–1,500 bp as background, SSP takes the average over a range of 500

98  kbp to 1 Mbp because we observed that the Jaccard score still decreases up to 1 Mbp (Fig.

99  1, step 3). Along with NSC and RSC, SSP also calculates "background uniformity" (Bu),

100  which evaluates the uniformity of mapped read distribution in background regions (step 4).

101  Finally, SSP calculates a "fragment cluster score" (FCS), which estimates the cluster level of

102  forward-reverse read pairs with each distance $d$ (step 5). FCS is the maximum difference in

103  the parameter cPNF (the cumulative proportion of neighboring fragments) at distance $d$

104    compared at background length. The outputs of SSP are displayed in PDF format and also

105    written to text files.

106

107    **Comparison with current methods**

108    To assess the performance of SSP for estimating S/N, we implemented three existing tools:

109    1) phantompeakqualtools (PPQT, https://github.com/kundajelab/phantompeakqualtools),

110    which internally implements spp version 1.14 [11] for cross-correlation analysis and then

111    outputs NSC and RSC; 2) Q version 1.2.0 [10], which adopts a strand-shift profile based on

112    the Hamming distance and calculates RSC; and 3) DeepTools version 2.5.0 [12], which

113    computes the synthetic Jensen-Shannon distance (JSD) that evaluates differences in the

114    cumulative fraction of mapped reads between ChIP by assuming a Poisson distribution as a

115    background model for windows of fixed length. We applied DeepTools with the "–

116    ignoreDuplicates" option according to the instructions given in the manual. We used default

117    parameters for each of the other tools.

118

119    **Estimating fragment length**

120    We first evaluated the performance of fragment-length estimation with SSP, PPQT and Q

121    using 65 paired-end ChIP-seq datasets for human, mouse, chicken, and fly (Fig. 2A and

122    Supplemental Table S1). We found that SSP could provide comparable and relatively more

123    accurate fragment-length data than PPQT and Q for all four species investigated. PPQT and

124    Q were nearly as accurate as SSP but could not provide a fragment-length estimate for

125    several of the samples (e.g., samples 37 and 45). On the other hand, none of the programs

126    could estimate an accurate fragment-length for certain samples (e.g., sample 16) for which

127    there was no clear peak in the strand-shift profile (Fig. 2B). Because it has been reported

128    that a high score for read-length shift can be mitigated by removing reads mapped on

5

129   "blacklist regions" in the genome [8], we re-analyzed 45 human samples (no. 1–45) after

130   removing reads mapped on blacklist regions [3] to validate the possibility that they affect the

131   accuracy of fragment-length estimation. However, such filtering had little effect

132   (Supplemental Fig. S2 and Supplemental Table S1). In fact, because the failure of fragment-

133   length estimation is mainly due to a lack of enrichment at the fragment-length shift, mitigating

134   the enrichment at read-length alone is insufficient. In this case, fragment length should be

135   supplied by the users. In subsequent analyses, we did not remove blacklist regions because

136   doing so could affect the RSC, and in fact detailed blacklist regions are available only for

137   human genome build hg19.

138

139   **Calculating the S/N for point and broad histone marks**

140   Required features for good S/N metrics are the quantifiability and sensitivity of different S/Ns

141   for both point- and broad-source factors, as well as the applicability to various cell types.  To

142   comprehensively evaluate the performance of SSP relative to other tools, we first used a

143   compendium of 860 ChIP-seq samples of histone modifications for 127 cell types, which

144   were obtained from the ROADMAP project [4]. These data contain information for six core

145   histone modifications, consisting of both point-source (H3K27ac, H3K4me1, H3K4me3) and

146   broad-source factors (H3K27me3, H3K36me3, H3K9me3) along with input samples. In the

147   consolidated dataset, reads of each sample were truncated to 36 bp, mapped onto genome

148   build hg19, filtered using a 36-bp mappability track, and then uniformly down-sampled to a

149   maximum depth of 30 million reads, which is appropriate for avoiding the effect derived from

150   different sequencing depths, parameters for mapping, and mappability.

151   A comparison is shown in Fig. 2C (see Supplemental Table S2 for detailed information and

152   scores for each sample). The results revealed that SSP-NSC and JSD could achieve

153   sufficient sensitivity both for point- and broad-source marks. The smaller difference between

154   point- and broad-source marks for JSD compared with SSP-NSC is perhaps a consequence

155    of score saturation, i.e., given that the maximum value of JSD is 1.0. PPQT-NSC showed

156    little difference among three broad marks compared with input samples (~1.1 fold), indicative

157    of insensitivity for broad marks.


158    As previously reported, RSCs obtained with all three tools were comparable or lower for

159    H3K9me3 than input samples. The discrepancy between NSC and RSC is possibly because

160    H3K9me3 is more highly enriched at the read-length shift compared with other histone

161    modifications derived from repetitive regions, such as centromeres [13]. Because RSC

162    amalgamates the magnitude of true peak enrichment and repeat effects, when the read-shift

163    enrichment is high, the RSC may be small even when the S/N is sufficiently high.

164    Furthermore, the relatively wider distribution of RSC for input samples indicates that a low

165    S/N increases the variability of it owing to the small value of the denominator (difference

166    between read-length value and background).


167    To further validate the ability of S/N indicators, we generated virtual data for histone

168    modifications with various S/Ns by adding a fixed number of input reads to each ChIP

169    sample in a stepwise manner. The S/N then decreased with increasing numbers of input

170    reads. Fig. 2D shows the comparison for E072 (Brain inferior temporal lobe) and

171    Supplemental Fig. S3 shows results for two other cell types. In most cases, the values of the

172    indicators decreased with increasing numbers of input reads. RSC was relatively higher for

173    H3K9me3 because, for this mark, the scores were often lower than those of the input (Fig.

174    2C). SSP-NSC had the superior or comparable sensitivity to changes in S/N, while PPQT-

175    NSC lacked sensitivity for evaluating broad marks.

176

177    **Evaluating the validity of the S/N for TFs for 20 cell types**

178    The S/N estimation could be affected by multiple factors, such as sequencing depth, read

179    length and copy number variations in cancer cell lines [14]. To validate the robustness of the

180    S/N indicators against these factors, we next investigated 399 ChIP-seq samples of TFs

181    (point source) for 20 cell types obtained from the ENCODE project [15]. This dataset

182    contains various read lengths (25, 36, and 50) and sequencing depths. Fig. 3A and

183    Supplemental Fig. S4 depict the distribution of SSP-NSC and the other scores, respectively,

184    for ChIP and input samples of 20 cell types (see Supplemental Table S3 for detailed

185    information). Whereas the number of samples varied among those cell types, we found that

186    SSP-NSC could reveal distinct differences between ChIP and input samples for all cell types.

187    To compare the various tools in this respect, we displayed the median scores for each cell

188    type for all indicators (Fig. 3B). For SSP-NSC and PPQT-NSC, median values for ChIP and

189    input samples were consistently different among all cell types, indicating that a cell type–

190    independent threshold value could be defined for these indicators. For example, SSP-NSC ≥

191    3.0 may be a good candidate threshold for TF ChIP samples, whereas the averaged S/N

192    varied among the TFs and antibodies used (Supplemental Fig. S4). Meanwhile, RSC and

193    JSD could not sufficiently distinguish ChIP and input samples. Although ChIP samples had

194    larger values than input samples for each cell type (Supplemental Fig. S5), the separation

195    between the data for ChIP and input samples depended on cell type, and therefore it was

196    difficult to determine a uniform threshold value. Consequently, SSP-NSC is a sensitive and

197    robust estimator that can be standardized across diverse cell types.

198

199    **Correlation with FRiP score**

200    To further evaluate the performance of S/N indicators, we calculated the Spearman's

201    correlation coefficient between the FRiP score and each S/N indicator across the ENCODE

202    and ROADMAP datasets (Table 1). Because FRiP score depends on sequencing depth, we

203    computed each FRiP score with and without total read normalization (see Methods for

204    details). First, RSC yielded a low correlation, suggesting that RSC cannot be used for

205    quantitative estimation of the S/N. In contrast, the output of each of SSP-NSC, PPQT-NSC,

206  and JSD was highly correlated with FRiP scores. Although SSP-NSC and PPQT-NSC each

207  correlated well with normalized FRiP scores, JSD correlated better with FRiP score without

208  normalization, which clearly shows the dependency of JSD on sequencing depth. This

209  conclusion is valid for the ROADMAP dataset with both point-source marks (H3K4me1,

210  H3K4me3, H3K27ac) and broad-source marks (H3K27me3, H3K36me3, H3K9me3) (Table

211  1). The lesser correlation of PPQT-NSC with broad-source marks compared with point-

212  source marks implies its lower sensitivity for broad-source marks.

213  To further investigate this tendency, we implemented a down-sampling analysis. We

214  selected six samples (four ChIP and two input samples) that contained an abundant number

215  of reads (>50 million) after removing duplicate reads. For each sample, we subsampled the

216  reads to a fixed number (from 5 million to 50 million) and calculated the ratio of the score at

217  each depth relative to the score for the 50 million reads (Fig. 3C and Supplemental Fig. S6A).

218  While all indicators except for JSD did not fluctuate with sequencing depth, JSD decreased

219  at lower sequencing depth. For input samples, each ratio fluctuated slightly (~1.1 fold)

220  because of smaller values for the 50 million reads. The analysis of histone modification data

221  also reached the same conclusion (Supplemental Fig. S6B). Consequently, SSP-NSC is the

222  best predictor of S/N for both point-source and broad-source marks, independent of

223  sequencing depth and cell types.

224

225  **Background uniformity**

226  NSC is defined as relative enrichment of the Jaccard score at each fragment-length shift

227  compared with the background level (Fig. 1, step 3). The next question was thus "why does

228  background level vary among samples?" By definition, the Jaccard score at background

229  reflects the co-occurrence probability of forward and reverse reads. Ideally, the background

230  reads should be uniformly distributed; in reality, however, the read distribution is often more

231  congregated, or biased, owing to various potential technical or biological issues [16],

9

232    resulting in a higher Jaccard score at background. Although the library complexity evaluates

233    the percentage of duplicate reads, it does not directly reflect any potential bias in the read

234    distribution. In fact, we observed that the background score increased ~2-fold when the

235    mapped reads were removed in every other 10-Mbp window, whereas library complexity and

236    NSC score remained essentially unchanged (Fig. 4A).

237    Based on this observation, we defined Bu, which evaluates the magnitude of the observed

238    background score compared with the uniform distribution (see Methods). A high value of Bu

239    indicates that the background reads are uniformly distributed even if library complexity is low.

240    In contrast, a low Bu score indicates sparse (or biased) read distribution, which decreases

241    the reliability of the peaks obtained.

242    We computed Bu scores for 860 histone modification samples from ROADMAP (Fig. 4B and

243    Supplemental Table S2). Although most of these consolidated data had library complexity =

244    1.0, we noted that a small amount of data for each histone modification and input sample

245    had a low Bu score (<0.8). A low Bu score was still observed even after filtering out samples

246    of low sequencing depth (<20 million reads). To further investigate the various aspects of Bu,

247    we chose 12 H3K36me3 samples as representatives, and the results are shown in Fig. 4C–

248    E. We grouped these samples into four types: (1) low NSC and high Bu, (2) high NSC and

249    high Bu, (3) high NSC and low Bu; this type was further classified as 3-1 (GC-rich) and 3-2

250    (not GC-rich). Fig. 4C illustrates the relative scores as a heatmap (see Supplemental Table

251    S4 for details concerning scores). Fig. 4D presents data for the read distribution proximal to

252    the housekeeping gene *IREB2* [17]. Groups 2 and 3 had high S/Ns, reflecting read

253    enrichment at the *IREB2* locus. Samples in group 3-1, however, had an unexpectedly sparse

254    read distribution, which is not reasonable considering that H3K36me3   is broadly distributed

255    within genic regions. Considering the GC-rich read distribution, this read distribution may be

256    a consequence of GC bias [18]. In contrast, group 3-2 had low Bu values without GC bias,

257    and read distribution was reasonable compared with group 3-1. However, this group also

258     had lower genome coverage in background region (Fig. 4E and Supplemental Fig. S7). A

259     possible reason for this is that the DNA fragmentation of tightly packed regions, e.g.,

260     heterochromatin, did not work well, resulting in a much lower number of reads on the regions.

261     These samples might confound the read normalization for comparative analyses that

262     assumes comparable read depth among samples over the entire genome [19]. These results

263     suggest that Bu is an effective criterion with which to judge whether a specific consideration

264     is required for comparative analysis.

265     Interestingly, GC-biased samples (group 3-1) had a striking peak for fragment length in the

266     strand-shift profile (Fig. 4F). This phenomenon might also facilitate the identification of read

267     bias.

268

269     **Relevance of Bu to other metrics**

270     To ascertain whether Bu varies among other mapping statistics and cell types, we next

271     investigated 399 ENCODE TF samples (Supplemental Table S3). We first found relatively

272     lower Bu values for MCF-7 cells (~0.8, Supplemental Fig. S8A), possibly owing to extensive

273     copy number variations [20]. The low-Bu samples also were more common when the S/N

274     was extremely high (e.g., RNA pol2, Supplemental Fig. S8B). Thus, it is desirable to use a

275     relaxed threshold value for Bu for these samples.

276     We next found that Bu did not correlate strongly with library complexity (Fig. 4G) or with the

277     mapping ratio of uniquely and multiply mapped reads (Supplemental Fig. S9). This result

278     suggests that the low values of these mapping statistics do not necessarily indicate biased

279     read distribution. For example, sample GM12892_PAX5-C20_v041610.1 had relatively low

280     library complexity (0.726) but a high Bu value (1.060) and no GC bias (peak = 45). The

281     strand-shift profile of this sample clearly revealed a maximum at fragment length

282     (Supplemental Fig. S10), indicative of sufficient quality.

283    On the other hand, Bu showed a moderately negative correlation with GC content (Fig. 4H),

284    consistent with the H3K36me3 results (Fig. 4C-E), whereas several samples had a high Bu

285    despite a highly GC-rich distribution (GC peak > 55); for instance, Rad21 sample for K562

286    (K562_Rad21_v041610.2) has GC peak = 56 but has an acceptable Bu (0.997). Although

287    Rad21 binding is closely correlated with CTCF [21], CTCF sample for K562

288    (K562_CTCF_SC-5916_PCR1x) is not GC-rich (GC peak = 48) and had a similar Bu (0.980).

289    In fact, this Rad21 sample had an unexpected bimodal GC distribution (Fig. 4I). Considering

290    the remarkable peak overlap between these two samples (98.6%, Supplemental Fig. S11),

291    the peaks of this Rad21 sample could be considered usable. This result implies that GC

292    content alone is not always appropriate to reject a putative low-quality sample. In this

293    respect, the Bu metric along with GC content provides a more reliable indicator of sample

294    quality with respect to biased read distribution.

295

296    **FCS can identify peak intensity and peak mode**

297    While having verified the effectiveness of SSP-NSC for calculating the S/N, we also found

298    that strand-shift profiles of a small number of input samples had peaks at fragment length

299    despite having a low FRiP score (e.g., input of E024 and E058 cells, Fig. 5A). These two

300    samples in particular had extremely high SSP-RSC (6.656 and 5.347), a phenomenon that is

301    commonly observed in PPQT and Q (Supplemental Fig. S12). We presumed that this is due

302    to "hidden duplicate reads". That is, at most two reads (forward and reverse pair) that are

303    derived from the same amplified DNA fragment can remain after PCR-bias filtering because

304    forward and reverse strands are scanned separately for single-end reads (Fig. 5B). Such

305    reads may often appear in low-library complexity samples and introduce a spike at the

306    fragment length, resulting in aberrant NSC and RSC values. To examine this hypothesis, we

307    generated strand-shift profiles for a paired-end sample in which both forward and reverse

308    reads were mapped as 'single-end'. As expected, the resulting profile showed a remarkable

309    peak at the fragment length shift (Fig. 5C). While NSC increased less drastically (1.53 to

310    2.54), RSC increased more than three times (0.61 to 2.29). This result suggested the

311    presence of the artifactual S/N enrichment without real peaks in a strand-shift profile, which

312    could especially influence the calculation of RSC.

313    To overcome this problem, we defined FCS, which directly evaluates the cluster level of

314    forward-reverse read pairs at each strand shift $d$ (see Methods for details). The FCS value is

315    high when read pairs with distance $d$ are highly clustered as peaks (Fig. 1, step 5). Therefore,

316    samples that contain hidden duplicate reads which are not clustered in a genome should

317    have a low FCS score. As expected, FCS could identify read clustering in samples and was

318    little affected by hidden duplicate reads (Fig. 5D). FCS correlated better with peak intensity

319    (height) than did FRiP, which represents a composite of peak number and intensity

320    (Supplemental Fig. S13).

321    Fig. 5E illustrates the example of five input samples from ROADMAP (see Supplemental

322    Table S5 for details concerning scores). The E097 input sample had strong peaks and the

323    highest FCS score among these samples (0.240). E024 and E058 (shown in Fig. 5A) had

324    high NSC and RSC values without many peaks, resulting in a low FCS score (0.041 and

325    0.038, respectively). In contrast, E100 had more peaks (33,476) than E097, but the FCS

326    score was low (0.044), indicating that the mapped reads were not highly clustered. The read

327    distribution and relatively lower FRiP score for E100 suggested that this sample had only

328    small peaks. Therefore, at a sufficiently high peak-calling threshold, most of the small peaks

329    (i.e., as in E100) would be expected to disappear, in contrast to the expectation for E097.

330    JSD was only minimally affected by hidden duplicate reads because it is not based on a

331    strand-shift profile, while it provided E100 with the highest score, suggesting that it

332    correlated better with peak number than did peak intensity and FRiP.

333    Interestingly, the FCS profile reflects the peak mode (point or broad source) for histone

334    modifications (Fig. 5F). H3K4me3 had the highest FCS at $d$ = fragment length and

335    decreased steeply at $d > 10$ kbp. The broad-source marks H3K27me3, H3K36me3, and

336    H3K9me3 each had a moderate score at fragment length, and the value was retained even

337    at $d > 10$ kbp, resulting in a higher score than for H3K4me3 at $d = 10$ kbp. H3K27ac had a

338    high score at fragment length and also the highest score at 10 kbp. This is not surprising

339    because H3K27ac had high peaks for point-source marks, some of which clustered in broad

340    genomic regions called super-enhancers [22]. This result suggested that FCS has the

341    potential to identify peak mode without the need for peak calling.

342

343    **Discussion**

344    The quality of ChIP-seq data depend on various experimental factors such as antibody

345    quality, crosslinking, DNA fragmentation, and PCR amplification. Although normalization

346    using a corresponding input sample mitigates biases in a ChIP sample, input data alone

347    cannot explain all the variability in read bias in the background [23]. It is important to assess

348    the genome-wide properties of samples in an objective manner to validate whether each

349    sample in the dataset requires special normalization or should be rejected for comparative

350    ChIP-seq analysis.

351    In this work, we present SSP, a peak calling–free quality assessment tool for read

352    enrichment in ChIP-seq data. We compared SSP against the existing methods PPQT, Q,

353    and DeepTools with more than 1,000 ChIP samples in public databases and demonstrated

354    that SSP has advantages over these methods with respect to sensitivity for both point-

355    source and broad-source factors, correlation with normalized FRiP, and robustness for

356    various sequencing depth and cell types. Although JSD, as utilized in DeepTools, is also

357    sensitive and can estimate the S/N for broad marks, it has less classification power between

358    ChIP and input samples owing to a lack of dynamic range. Moreover, because JSD depends

359    on sequencing depth, it requires subsampling for comparison across samples, which is

360    burdensome for a large-scale analysis.

361  Bu evaluates the reliability of the obtained peaks by quantifying the distribution of mapped

362  reads in background regions. Although GC content correlates with the bias level in ChIP

363  samples, it alone cannot be used for filtering because samples that have many GC-rich

364  peaks (e.g., CpG islands) also have a high GC content. Bu is beneficial in this regard,

365  especially for consolidated data, for which the mapping ratio and library complexity metrics

366  are not generally available. While the "X-intercept" metric in DeepTools evaluates genome

367  coverage, it also depends on sequencing depth and less robust than Bu. Finally, SSP

368  provides FCS, which avoids the effect of hidden duplicate reads. The potential of FCS to

369  evaluate peak mode may facilitate capturing dynamic changes of genome-wide binding

370  patterns among samples, such as during the cell development [24].

371  Owing to the difficulty of assessing broad marks and peak reliability, a previous study

372  involving large-scale sample evaluation for S/N was limited to input and negative-control

373  samples [25]. The use of SSP enables in-depth validation using >1,000 ChIP-seq data that

374  are publicly available, including point-source and broad-source marks, along with virtual data,

375  and SSP provides multiple key insights for ChIP-seq analysis.

376  Based on our results, we recommend using NSC rather than RSC when calculating the S/N

377  in the strand-shift profile for several reasons. First, RSC is based on the value at read length,

378  which depends on blacklist region filtering. Second, RSC has high variance in the evaluation

379  of low-S/N samples due to the small values at both read-length and fragment-length. Third,

380  RSC combines the magnitude of peak enrichment and repeat effects. A strong repeat effect

381  cancels out strong peak intensity. Finally, we observed that, compared with NSC, RSC is

382  strongly affected by hidden duplicate reads.

383  One challenge that remains is to identify false-positive peaks caused by non-specific binding,

384  such as "hyper-ChIPable regions" [26]. SSP and all existing tools cannot distinguish whether

385  or not DNA-binding is derived from true binding, and thus a comparison with mock ChIP-seq

15

386    data (e.g., IgG) is needed to avoid such false positives. Finally, the challenge remains to

387    accurately estimate fragment length from single-end data.

388

389    **Methods**

390    **Strand-shift profile using the Jaccard index**

391    Let $v_c^{fwd} = \left(x_1^{fwd} x_2^{fwd} \dots x_n^{fwd}\right)_c$ and $v_c^{rev} = (x_1^{rev} x_2^{rev} \dots x_n^{rev})_c$ be strand-specific binary

392    vectors for forward and reverse strands for chromosome *c* of length *n*, respectively. $x_k^{str}$

393    ($k \in [1, n]$, $str \in [fwd, rev]$) is the number of reads whose 5' ends map to position *k* of strand

394    *str*, and $x_k^{str} \in \{0,1\}$ after removing duplicate reads. The Jaccard index between $v_c^{fwd}$ and

395    $v_c^{rev}$ at strand shift *d* is defined as follows:

396    $$J[v^{fwd}, v^{rev}, d]_c = \frac{|\, v_c^{fwd} \cap v_c^{rev}(d)|}{|\, v_c^{fwd} \cup v_c^{rev}(d))|},$$

397    where $v_c^{rev}(d)$ is $(x_{d+1}^{rev} x_{d+2}^{rev} \dots x_n^{rev})_c$. Therefore, $0 \le J[v^{fwd}, v^{rev}, d]_c \le 1$. This formula can

398    be transformed as follows:

399    $$J[v^{fwd}, v^{rev}, d]_c = N_d^{both}/(N^{fwd} + N^{rev} - N_d^{both}),$$

400    where $N^{str} = \sum_k^n x_k^{strand}$, and $N_d^{both} = \sum_k^{n-d} f\left(x_k^{fwd}, x_{k+d}^{rev}\right)$, where $f(a, b) = \begin{cases} 1 \ (a = b = 1) \\ 0 \ (\text{otherwise}) \end{cases}$. This

401    score is calculated using the bitset operator in C++. The strand shift *d* ranges from –500 bp

402    to 1,500 bp at single–base pair resolution. To standardize the value for various species

403    having different genome lengths, this Jaccard score is then normalized per fixed number of

404    reads ($N_{const}$, 10M default) for a fixed length of bases ($L_{const}$, 100M default):

405    $$Jnorm[v^{fwd}, v^{rev}, d]_c = J[v^{fwd}, v^{rev}, d]_c * \frac{N_{const}}{N_c} * \frac{L_{const}}{L_c},$$

16

406    where $N_c$ and $L_c$ are the number of mapped reads and the number of mappable positions (at

407    which the reads starting at those positions are uniquely mapped on the genome),

408    respectively, on chromosome *c.* We estimated $L_c$ for 36-mer and 50-mer reads based on the

409    code from Peakseq [27].

410    Finally, SSP assembles the Jaccard index profiles obtained from all autosomes:

411
$$Jnorm[v^{fwd}, v^{rev}, d]_{genome} = \sum_{c \in C} \frac{N_c}{N_{genome}} Jnorm[v^{fwd}, v^{rev}, d]_c,$$

412    where *C* is the set of all autosomes, and $N_{genome} = \sum_{c \in C} N^c$. SSP excludes sex

413    chromosomes to ignore gender-specific differences. We use this $Jnorm[v^{fwd}, v^{rev}, d]_{genome}$

414    as the Jaccard score $J(d)$ for each sample in SSP. Then the fragment length $d_{flen}$ can be

415    estimated as $d_{flen} = argmax_{d_{readlen}*1.2 < d < 1500} J(d)$. To ignore a peak at the read-length shift

416    ($d_{readlen}$), SSP uses $d > d_{readlen} * 1.2$. Then NSC and RSC can be calculated as:

417
$$NSC = J(d_{flen})/J(d_{bg}) \text{ and } RSC = \left(J(d_{flen}) - J(d_{bg})\right) / \left(J(d_{readlen}) - J(d_{bg})\right),$$

418    where $J(d_{bg})$ is the Jaccard score for the background, which is the average from 500 kbp to

419    1 Mbp at steps of 5 kbp (default).

420

421    **Background uniformity**

422    Background uniformity *Bu* is defined as follows:

423
$$Bu = J(d_{bg})^{uniform} / J(d_{bg})^{observe},$$

424    where $J(d_{bg})^{uniform}$ is the normalized Jaccard score for the background for a sample that

425    has a completely uniform read distribution. That is, by denoting $E[x_k = 1]^{strand}$ as the

426    probability of a mapped read occurring at genomic position *k*,

$$J(d_{bg})^{uniform} = \frac{\left(E[x_k = 1]^{fwd} * E[x_k = 1]^{rev} * L_c\right)}{\left(\frac{N_{const}}{2} + \frac{N_{const}}{2} - E[x_k = 1]^{fwd} * E[x_k = 1]^{rev} * L_{const}\right)}$$

427 $$= N_{const}/(4L_{const} - N_{const}),$$

428 because $E[x_k = 1]^{fwd} = E[x_k = 1]^{rev} = \frac{N_{const}}{2L_{const}}$. Therefore, $J(d_{bg})^{uniform} = 1/39$ for $N_{const} = $

429 10M and $L_{const}$ = 100M. A high *Bu* score indicates that the sample has a relatively uniform

430 read distribution in the background region. *Bu* should range from 0 to 1, but practically, the

431 maximum score of Bu slightly exceeds 1.0 because the estimated mappable chromosomal

432 length $L_c$ is a bit larger than the actual mappable chromosomal length.

433

434 **FCS**

435 Similar to the Jaccard score, FCS is calculated for each strand shift *d*. Let $rp(d)$ represent a

436 forward- and reverse-read pair with distance *d*. Denote all read-pair sets as

437 $\left\{rp_1(d), rp_2(d), \dots, rp_{N_d^{both}}(d)\right\}$, which is sorted by genomic position. That is, this set consists

438 of all $(x_k^{fwd}, x_{k+d}^{rev})$ pairs such that $f\left(x_k^{fwd}, x_{k+d}^{rev}\right) = 1$.

439 Let $NF[d, s]$ represent the number of $rp_k(d)$ ($k \in \left[1, N_d^{both} - 1\right]$) that have neighboring read

440 pairs $rp_{k+1}(d)$ within distance *s*. Then the cPNF is:

$$cPNF[d, s] = NF[d, s]/NF[d, L_c]$$

441 $$= NF[d, s]/N_d^{both}.$$

442 This cPNF is calculated up to $s_{max}$ (5 kbp, default). Supplemental Fig. 14A shows the typical

443 pattern of cPNF for ChIP and input samples. If a sample has peaks, the cPNF score

444 becomes higher at short distance *d* (Supplemental Fig. 14A left). If a sample does not have

445    peaks, the cPNF score at short distance shows little difference from that at long distance

446    (background) (Supplemental Fig. 14A right).

447    Then, we define FCS[d] as the maximum difference of cPNF[d] from background against *s*:

448    $$\mathrm{FCS[d]} = \arg\max_{0 \leq s < s_{max}} \left( \mathrm{cPNF[d}, s] - \mathrm{cPNF}\left[\mathrm{d_{bg}}, s\right] \right).$$

449    Because FCS[d] depends on sequencing depth, SSP down-samples the reads to a fixed

450    number (10M reads, default). The maximum difference strategy used here can provide more

451    robust values for different values of parameter $s_{max}$ and $d_{bg}$ than a relative entropy method

452    such as the Kullback–Leibler divergence. The resulting FCS profile (Supplemental Fig. 14B)

453    reflects the cluster level of $rp(d)$ in the sample, whereas the Jaccard score $J(d)$ reflects

454    $N_d^{both}$, the number of $rp(d)$. Therefore, for point-source factors, the average peak width is ~1

455    kbp, and $\mathrm{FCS[d]}(d \leq 1\,\mathrm{kbp})$ is high, whereas broad marks have relatively higher scores for

456    the broad width (e.g., H3K4me3 and H3K36me3 in Fig. 5F).

457

458    **Read mapping**

459    We used bowtie 1 version 1.1.2 [28] for mapping single-end reads and extracted uniquely

460    mapped reads. For mapping paired-end reads, we used bowtie 2 version 2.2.9 [29], which is

461    more sensitive for longer reads than bowtie 1.

462

463    **FRiP score**

464    We used MACS2 version 2.1.1 [30] for peak calling with --nomodel option, and we also

465    supplied --broad option for broad marks (H3K27me3, H3K36me3, H3K9me3). Because

466    MACS2 does not have an option for read-depth normalization, we also used DROMPA3

467    version 3.2.6 [31] with "–n GR" option for peak calling normalized by the number of

468    nonredundant reads. FRiP scores with and without total read normalization were calculated

469    by peaks of MACS2 and DROMPA3, respectively. The FRiP score, peak height, library

470    complexity for 10M reads, and GC content were also calculated with DROMPA3.

471

472    **Data access**

473    For histone modification data, we acquired the consolidated data for 117 cell types (tagAlign

474    format, build hg19) from the ROADMAP project [32], available at

475    http://egg2.wustl.edu/roadmap/web_portal/. For the analysis of Q and DeepTools, we

476    converted tagAlign format to BAM format using bedtools (http://bedtools.readthedocs.io). For

477    TF data for the 20 cell types, we acquired fastq files from the Sequence Read Archive (SRA)

478    under accession number SRP008797, which is part of the ENCODE project [15].

479    Supplemental Methods describes the method for generating virtual data (Fig. 2D, Fig 3C and

480    Fig 4A).

481

482    **Software availability**

483    SSP is open-source software that is freely available for nonprofit use. It is implemented as a

484    C++ package with a boost library (http://www.boost.org/), and internally uses R to visualize

485    the strand-shift profile and FCS profiles in PDF format. The user manual and examples are

486    available at https://github.com/rnakato/SSP. The mappability tables generated for several

487    species are also provided on the SSP website.

488

489    **Competing interests**

490    The authors declare that they have no competing interests.

491

## Acknowledgements

495

## FUNDING

500

## REFERENCES

502 1.     Park PJ: ChIP-seq: advantages and challenges of a maturing technology. *Nat Rev Genet*
503     2009, 10:669-680.
504 2.     Furey TS: ChIP-seq and beyond: new and improved methodologies to detect and
505     characterize protein-DNA interactions. *Nat Rev Genet* 2012, 13:840-852.
506 3.     Encode Project Consortium: An integrated encyclopedia of DNA elements in the human
507     genome. *Nature* 2012, 489:57-74.
508 4.     Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,
509     Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, et al: Integrative analysis of 111
510     reference human epigenomes. *Nature* 2015, 518:317-330.
511 5.     Stunnenberg HG, International Human Epigenome C, Hirst M: The International Human
512     Epigenome Consortium: A Blueprint for Scientific Collaboration and Discovery. *Cell* 2016,
513     167:1145-1149.
514 6.     Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglou S, Bernstein BE, Bickel
515     P, Brown JB, Cayting P, et al: ChIP-seq guidelines and practices of the ENCODE and
516     modENCODE consortia. *Genome Res* 2012, 22:1813-1831.
517 7.     Bailey T, Krajewski P, Ladunga I, Lefebvre C, Li Q, Liu T, Madrigal P, Taslim C, Zhang J:
518     Practical guidelines for the comprehensive analysis of ChIP-seq data. *PLoS Comput Biol*
519     2013, 9:e1003326.
520 8.     Carroll TS, Liang Z, Salama R, Stark R, de Santiago I: Impact of artifact removal on ChIP
521     quality metrics in ChIP-seq and ChIP-exo data. *Front Genet* 2014, 5:75.
522 9.     Nakato R, Shirahige K: Recent advances in ChIP-seq analysis: from quality management to
523     whole-genome annotation. *Brief Bioinform* 2016.

524 10.    Hansen P, Hecht J, Ibrahim DM, Krannich A, Truss M, Robinson PN: Saturation analysis of
525        ChIP-seq data for reproducible identification of binding peaks. *Genome Res* 2015, 25:1391-
526        1400.
527 11.    Kharchenko PV, Tolstorukov MY, Park PJ: Design and analysis of ChIP-seq experiments for
528        DNA-binding proteins. *Nat Biotechnol* 2008, 26:1351-1359.
529 12.    Ramirez F, Ryan DP, Gruning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S, Dundar F,
530        Manke T: deepTools2: a next generation web server for deep-sequencing data analysis.
531        *Nucleic Acids Res* 2016, 44:W160-165.
532 13.    Guenatri M, Bailly D, Maison C, Almouzni G: Mouse centric and pericentric satellite
533        repeats form distinct functional heterochromatin. *J Cell Biol* 2004, 166:493-505.
534 14.    Zarrei M, MacDonald JR, Merico D, Scherer SW: A copy number variation map of the
535        human genome. *Nat Rev Genet* 2015, 16:172-183.
536 15.    Gertz J, Savic D, Varley KE, Partridge EC, Safi A, Jain P, Cooper GM, Reddy TE, Crawford GE,
537        Myers RM: Distinct properties of cell-type-specific and shared transcription factor binding
538        sites. *Mol Cell* 2013, 52:25-36.
539 16.    Meyer CA, Liu XS: Identifying and mitigating bias in next-generation sequencing methods
540        for chromatin biology. *Nat Rev Genet* 2014, 15:709-721.
541 17.    Eisenberg E, Levanon EY: Human housekeeping genes, revisited. *Trends Genet* 2013,
542        29:569-574.
543 18.    Benjamini Y, Speed TP: Summarizing and correcting the GC content bias in high-
544        throughput sequencing. *Nucleic Acids Res* 2012, 40:e72.
545 19.    Lun AT, Smyth GK: De novo detection of differentially bound regions for ChIP-seq data
546        using peaks and windows: controlling error rates correctly. *Nucleic Acids Res* 2014, 42:e95.
547 20.    Hampton OA, Den Hollander P, Miller CA, Delgado DA, Li J, Coarfa C, Harris RA, Richards S,
548        Scherer SE, Muzny DM, et al: A sequence-level map of chromosomal breakpoints in the
549        MCF-7 breast cancer cell line yields insights into the evolution of a cancer genome.
550        *Genome Res* 2009, 19:167-177.
551 21.    Wendt KS, Yoshida K, Itoh T, Bando M, Koch B, Schirghuber E, Tsutsumi S, Nagae G,
552        Ishihara K, Mishiro T, et al: Cohesin mediates transcriptional insulation by CCCTC-binding
553        factor. *Nature* 2008, 451:796-801.
554 22.    Hnisz D, Abraham BJ, Lee TI, Lau A, Saint-Andre V, Sigova AA, Hoke HA, Young RA: Super-
555        enhancers in the control of cell identity and disease. *Cell* 2013, 155:934-947.
556 23.    Rashid NU, Giresi PG, Ibrahim JG, Sun W, Lieb JD: ZINBA integrates local covariates with
557        DNA-seq data to identify broad and narrow regions of enrichment, even within amplified
558        genomic regions. *Genome Biol* 2011, 12:R67.
559 24.    Dahl JA, Jung I, Aanes H, Greggains GD, Manaf A, Lerdrup M, Li G, Kuan S, Li B, Lee AY, et al:
560        Broad histone H3K4me3 domains in mouse oocytes modulate maternal-to-zygotic
561        transition. *Nature* 2016, 537:548-552.
562 25.    Marinov GK, Kundaje A, Park PJ, Wold BJ: Large-Scale Quality Analysis of Published ChIP-
563        seq Data. *G3 (Bethesda)* 2014, 4:209-223.
564 26.    Teytelman L, Thurtle DM, Rine J, van Oudenaarden A: Highly expressed loci are vulnerable
565        to misleading ChIP localization of multiple unrelated proteins. *Proc Natl Acad Sci U S A*
566        2013, 110:18602-18607.
567 27.    Rozowsky J, Euskirchen G, Auerbach RK, Zhang ZD, Gibson T, Bjornson R, Carriero N,
568        Snyder M, Gerstein MB: PeakSeq enables systematic scoring of ChIP-seq experiments
569        relative to controls. *Nat Biotechnol* 2009, 27:66-75.
570 28.    Langmead B, Trapnell C, Pop M, Salzberg SL: Ultrafast and memory-efficient alignment of
571        short DNA sequences to the human genome. *Genome Biol* 2009, 10:R25.
572 29.    Langmead B, Salzberg SL: Fast gapped-read alignment with Bowtie 2. *Nat Methods* 2012,
573        9:357-359.

574    **30.**    Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM,
575         Brown M, Li W, Liu XS: Model-based analysis of ChIP-Seq (MACS). *Genome Biol* 2008,
576         9:R137.
577    **31.**    Nakato R, Itoh T, Shirahige K: DROMPA: easy-to-handle peak calling and visualization
578         software for the computational analysis and validation of ChIP-seq data. *Genes Cells* 2013,
579         18:589-601.
580    **32.**    Zhou X, Li D, Zhang B, Lowdon RF, Rockweiler NB, Sears RL, Madden PA, Smirnov I, Costello
581         JF, Wang T: Epigenomic annotation of genetic variants using the Roadmap Epigenome
582         Browser. *Nat Biotechnol* 2015, 33:345-346.

583

584

585

586

**Figure 1.** Workflow of SSP. Step 1: convert mapped reads to strand-specific binary vectors

($n$: chromosome length), in which '1' indicates that the 5' end of a read is mapped at the

genomic position. Duplicate reads are discarded. Step 2: calculate the similarity between

forward and reverse strand-specific binary vectors for each strand shift $d$ based on the

Jaccard index. An example calculation is shown ($n = 10$, $d = 0, 1, 2$). Step 3: plot a strand-

shift profile based on the Jaccard index and calculate NSC and RSC. Fragment length is

estimated as the distance $d$ at which the Jaccard score is maximal except for read-length

shift. Step 4: calculate background uniformity based on the background level. Step 5:

calculate the fragment cluster score to evaluate the cluster level of all forward-reverse read

pairs with each distance $d$ (orange rectangles). These read pairs are the same as the red

bars in step 2. The variable s indicates the distance to the nearest downstream read pair.

24

598     cPNF is the cumulative proportion of neighboring downstream fragments (see Methods).

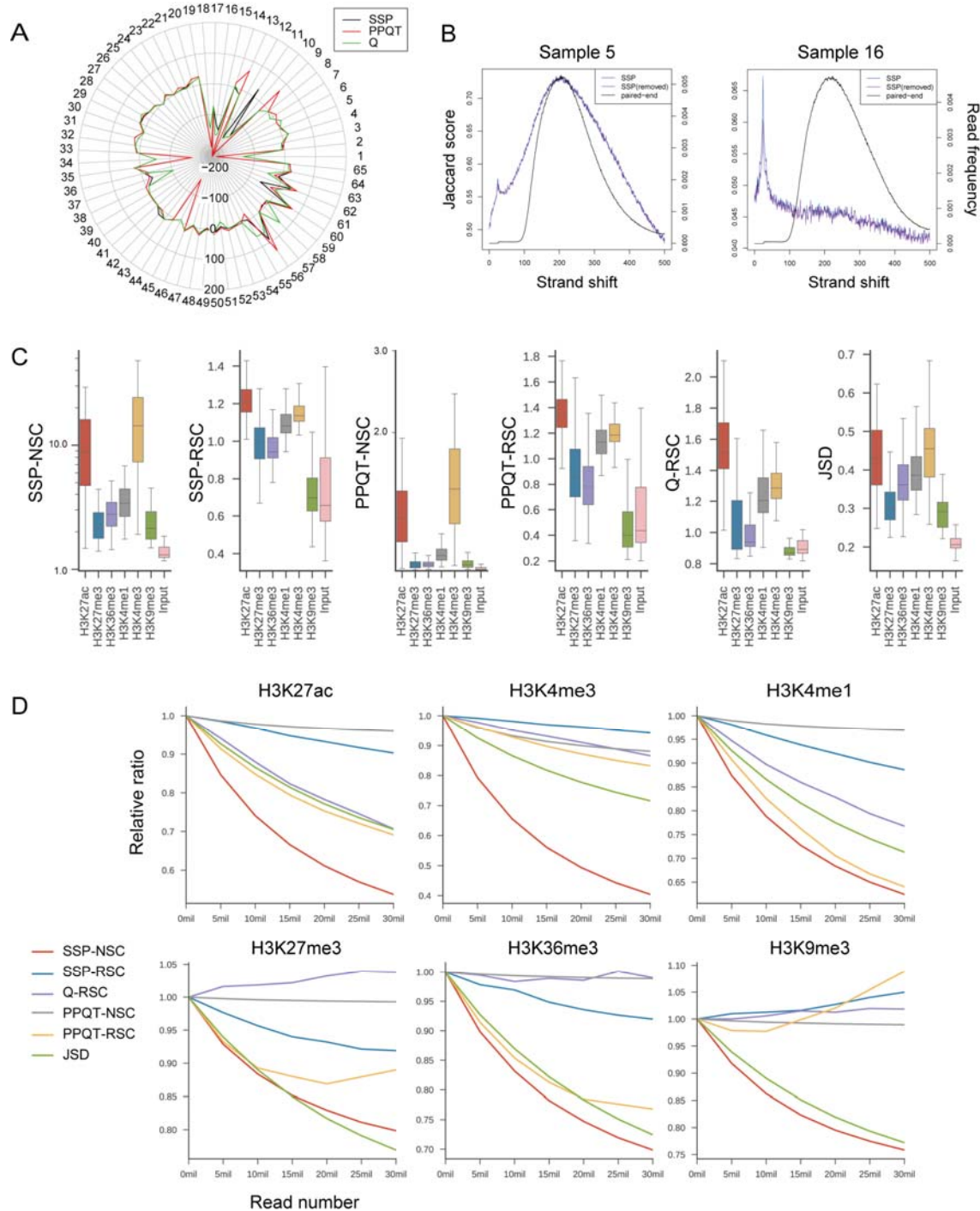599     FCS increases as read pairs become more clustered.

600

601

**Figure 2.** (A) Radar plot of the comparison between the fragment-length estimated by each tool and that from paired-end data for 65 paired-end ChIP-seq data for human (1–45, mapped to build hg38), mouse (46–54, mm10), chicken (55–61, galGal4), and fly (62–65, dm6). The y axis indicates the difference between the fragment size estimated from the

607     single-end (F3) reads by these tools and that derived from the paired-end reads. See

608     Supplemental Table S1 for names and scores for each sample. (B) Examples of strand-shift

609     profiles by SSP for original data (blue) and after removing blacklist regions (purple), with

610     fragment length distribution inferred by paired-end data (black). Samples that have a clear

611     peak (left) for fragment length can be estimated accurately. The estimation is less accurate if

612     there is no peak and much repetitive enrichment in the sample (right). (C) The distribution of

613     scores by SSP (NSC and RSC), PPQT (NSC and RSC), Q (RSC), and DeepTools (JSD).

614     Note that the y axis is a log-scale for SSP-NSC and PPQT-NSC but a linear scale for the

615     others. (D) Relative ratio at each different S/N (adding input reads from 5 million to 30

616     million) against original data. Data were obtained from cell type E072 (Brain inferior temporal
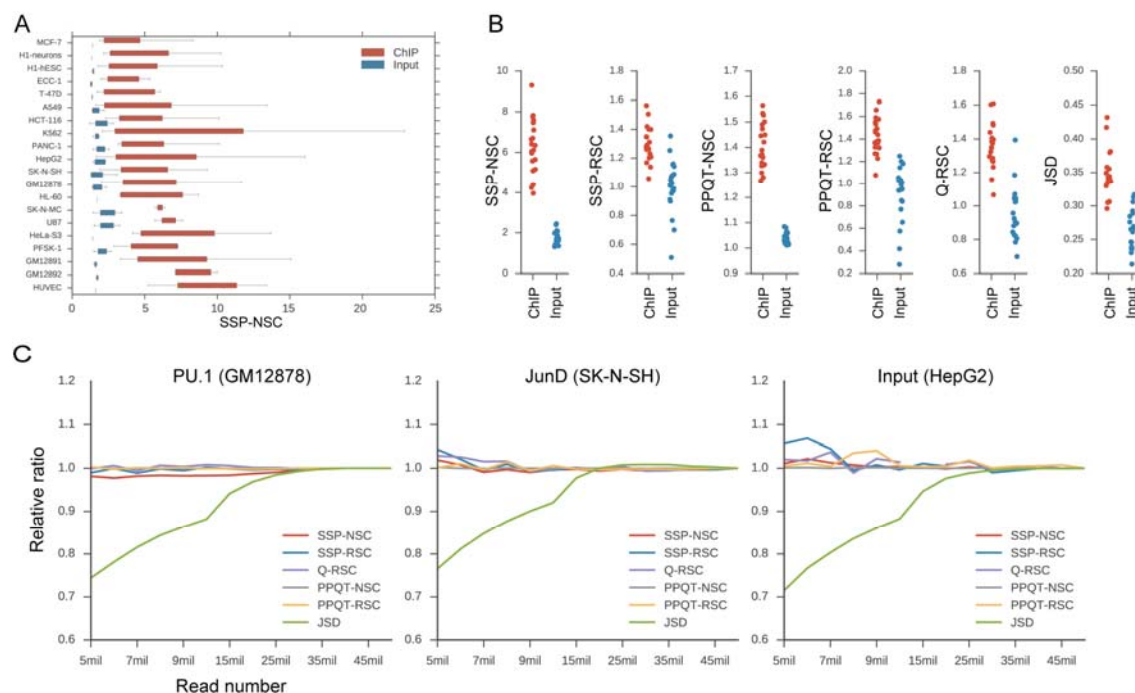
617     lobe) of ROADMAP.

618

**Figure 3.** (A) SSP-NSC distribution of ChIP (red) and input (blue) for 20 cell types. (B) Median values of S/N indicators for 20 cell types. (C) Relative S/Ns at each sequencing depth (5–50 million) against 50 million reads. Duplicate reads were removed in advance.
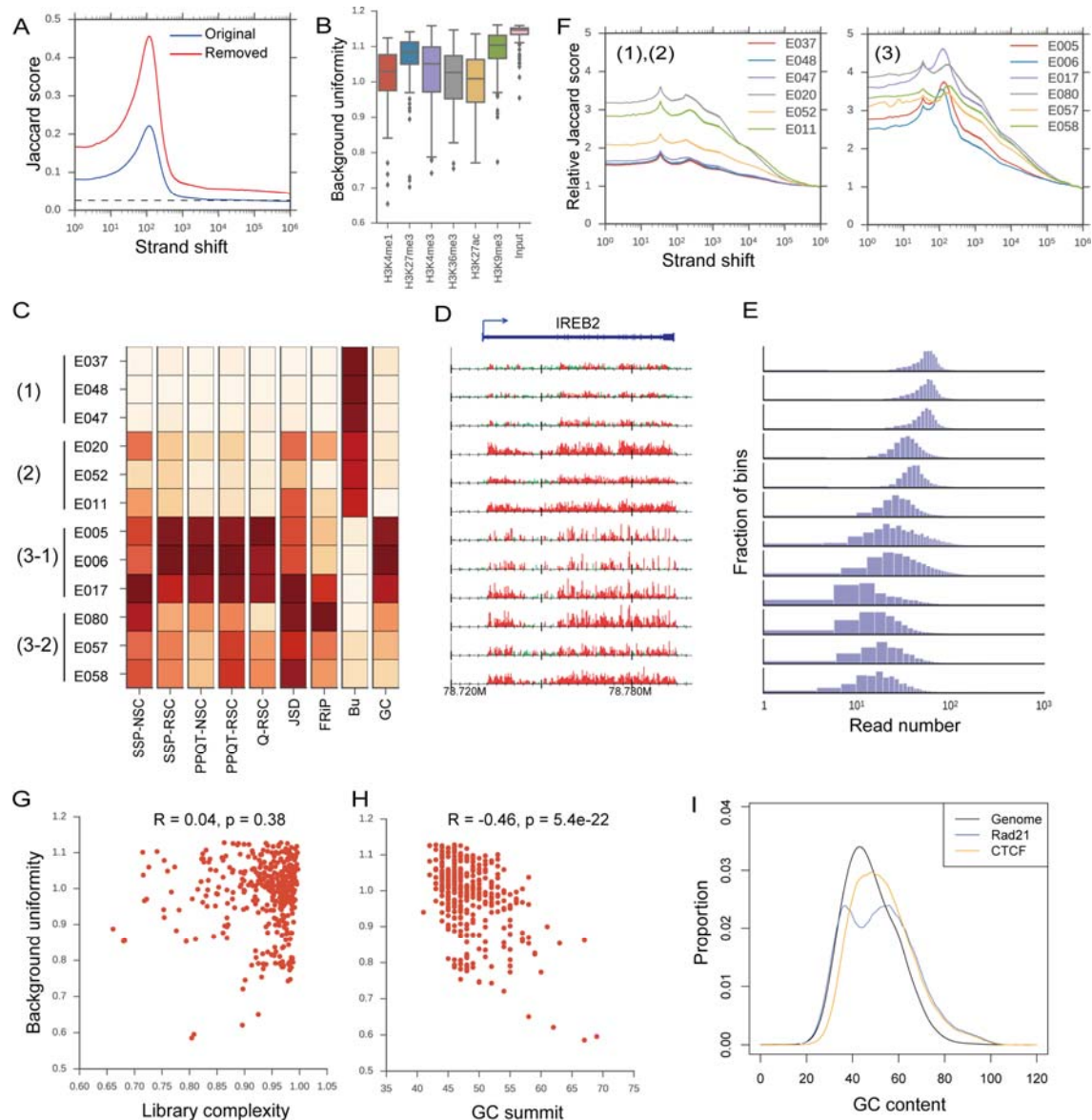
28

**Figure 4.** (A) Strand-shift profile for PU.1 data for GM12878 cells from ENCODE. Blue,

original data. Red, data after removing mapped reads in every other 10-Mbp window. The

horizontal dashed line indicates the expected background level. (B) Distribution of

background uniformity for histone modifications. (C–F) Analysis of H3K36me3 data for 12

cell types from ROADMAP. (C) Heatmap of S/N scores alongside Bu and GC peak. Darker

colors indicate higher values. See Supplemental Table S4 for the description and detailed

scores of each sample. (D) The read distribution around the IREB2 locus (chromosome 15,

78.72–78.80 Mbp). Read number was normalized by the total number of nonredundant

633    reads. The peak regions identified by MACS2 are highlighted in red. (E) Histogram of

634    mapped read number for each 100-kbp bin of the whole genome except chromosome Y. (F)

635    Strand-shift profiles for the relative Jaccard score against background for groups 1 and 2

636    (left) and group 3 (right). (G, H) Correlation plot between Bu and library complexity (G) and

637    GC peak (H). (I) Distribution of GC content over the entire genome (black),

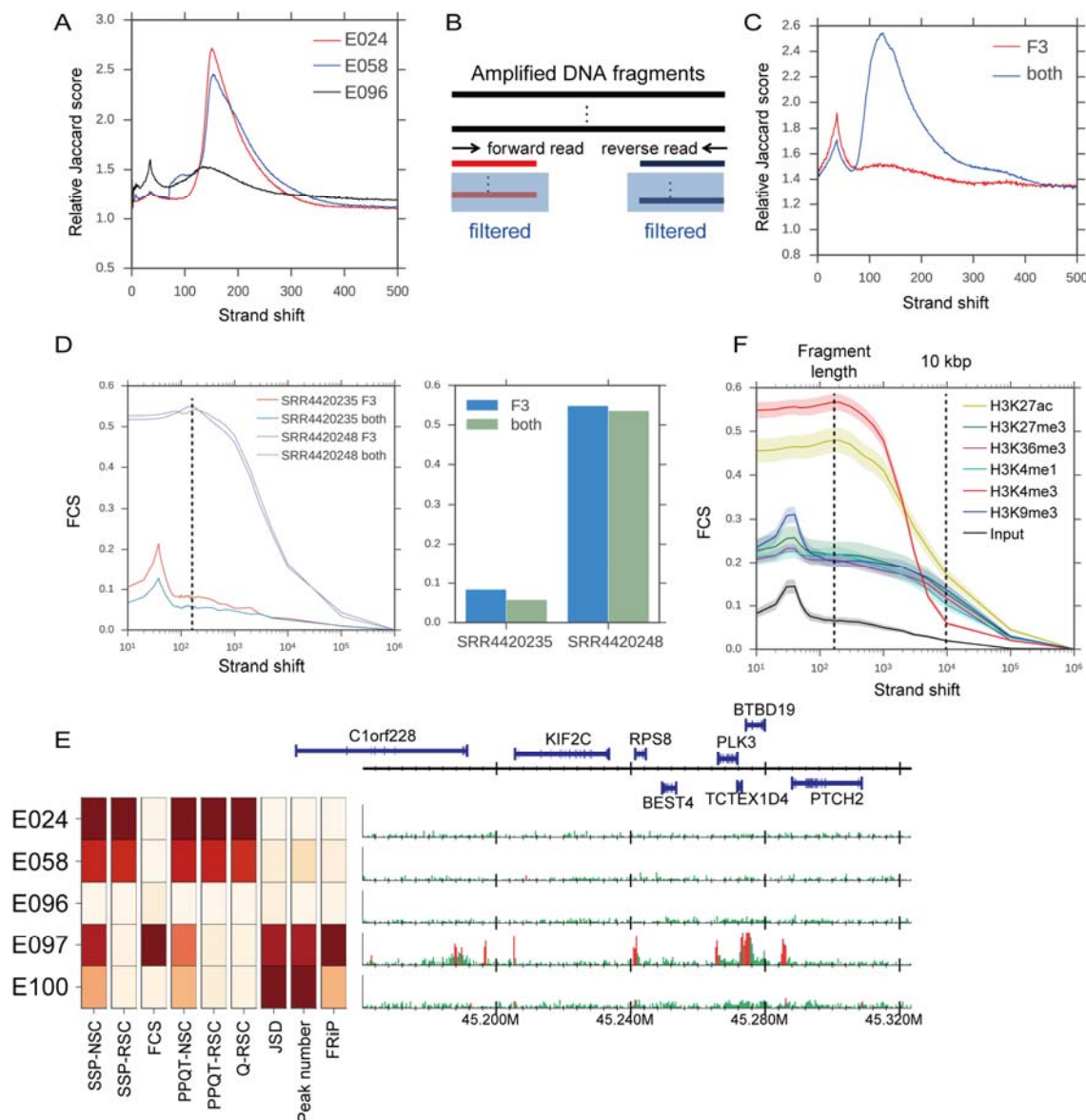638    Rad21_v041610.2 (blue), and CTCF_SC-5916_PCR1x (yellow) for K562 cells.

639

640

641

**Figure 5.** (A) Strand-shift profiles (0 < d < 500) for two input samples (E024, E058), which

have apparent peaks at fragment length but actually have low FRiP score (< 0.01). The y-

axis indicates the relative Jaccard score against the background level. A typical profile for

input (E096) that has a similar FRiP score is also shown. (B) Schematic illustration of hidden

duplicate reads. (C) Strand-shift profile for sample SRR4420235 (no. 25 in Fig. 2A) using F3

read only (red) and both F3 and F5 reads as single-end (blue). (D) FCS distribution for each

strand shift for each of samples SRR4420235 and SRR4420248 (no. 31 in Fig. 2A) (left).

The value at fragment length shown as a dashed vertical line is used as the FCS score

650    (right). (E) Heatmap for each S/N value and number of peaks identified by MACS2 of five

651    input samples, and their read distribution (chromosome 1, 45.16–45.32 Mbp). The peak

652    regions are highlighted in red. (F) FCS profile for all ROADMAP data. Lines and shaded

653    regions indicate the mean value and 95% confidence interval, respectively, at the each

654    strand shift (x-axis).

655

656

657    **Table 1.** Spearman's correlation between each S/N and FRiP score without read

658    normalization (top) and with normalization (bottom). *p-value for correlation coefficient < 0.01.

|  | SSP-NSC | SSP-RSC | PPQT-NSC | PPQT-RSC | Q-RSC | JSD | FCS |
|---|---|---|---|---|---|---|---|
| ENCODE | *0.84  *0.90 | 0.10  *0.28 | *0.83  *0.88 | *0.15  *0.31 | *0.16  *0.32 | *0.93  *0.81 | *0.79  *0.80 |
| ROADMAP (point source) | *0.77  *0.94 | *0.20  *0.28 | *0.70  *0.89 | *0.33  *0.34 | *0.29  *0.30 | *0.97  *0.90 | *0.44  *0.73 |
| ROADMAP (broad source) | *0.72  *0.89 | *0.40  *0.31 | *0.43  *0.64 | *0.37  *0.32 | *0.37  *0.31 | *0.90  *0.79 | *0.55  *0.88 |

659