1  **SMARTcleaner: identify and clean off-target signals in**

2  **SMART ChIP-seq analysis**

3

4  Dejian Zhao[1] and Deyou Zheng[1,2,3,*]

5

6  [1]Department of Genetics, [2]Department of Neurology, [3]Department of

7  Neuroscience, Albert Einstein College of Medicine, 1300 Morris Park Ave., Bronx,

8  New York, USA

9

10  D. Zhao: dejian.zhao@einstein.yu.edu

11  D. Zheng: deyou.zheng@einstein.yu.edu

12

13  *Correspondence should be addressed to,

14  Deyou Zheng, Ph.D.

15  Tel: +1 718 678 1217

16  Fax: +1 718 430 8785

17  Email: deyou.zheng@einstein.yu.edu

18

19

20

21

22

## Abstract

**Background**: Noises and artifacts may arise in several steps of the next-generation sequencing (NGS) process. Recently, a NGS library preparation method called SMART, or _S_witching _M_echanism _A_t the 5' end of the _R_NA _T_ranscript, is introduced to prepare ChIP-seq (chromatin immunoprecipitation and deep sequencing) libraries from small amount of DNA material. The protocol adds Ts to the 3' end of DNA templates, which is subsequently recognized and used by SMART poly(dA) primers for reverse transcription and then addition of PCR primers and sequencing adapters. The poly(dA) primers, however, can anneal to poly(T) sequences in a genome and amplify DNA fragments that are not enriched in the immunoprecipitated DNA templates. This off-target amplification results in false signals in the ChIP-seq data.

**Results**: Here, we show that the off-target ChIP-seq reads derived from false amplification of poly(T/A) genomic sequences have unique and strand-specific features. Accordingly, we develop a tool (called "SMARTcleaner") that can exploit the features to remove SMART ChIP-seq artifacts. Application of SMARTcleaner to several SMART ChIP-seq datasets demonstrates that it can remove reads from off-target amplification effectively, leading to improved ChIP-seq peaks and results.

**Conclusions**: SMARTcleaner could identify and clean the false signals in SMART-based ChIP-seq libraries, leading to improvement in peak calling, and downstream data analysis and interpretation.

46    **Keywords**:

47    SMART, ChIP-seq, NGS, false priming, false amplification

48

## 49    Background

50    In the past decade, deep sequencing by next generation sequencing (NGS) has

51    been widely applied in nearly all fields of biological research, in which information

52    from biological processes (e.g., transcription and protein-DNA interaction) can be

53    converted to DNAs for sequencing [1-4]. NGS is a complex procedure involving

54    DNA/RNA isolation, library preparation, deep sequencing, data processing and

55    interpretation. Each of these steps can introduce biases and artifacts, but the first

56    step - preparation of NGS libraries is arguably the most critical phase as errors

57    can be propagated to later steps, if not carefully controlled [5, 6]. Among them,

58    PCR amplification is a major source of bias due to the fact that not all fragments

59    are amplified with the same efficiency [5].

60

61    As powerful as NGS technology is, its application with limited amounts of

62    biological material, for example, DNA or RNA isolated from a very small number

63    of cells, remains a challenge. This is primarily due to the low efficiency in ligating

64    targeted DNA/RNA fragments to the NGS sequencing adaptors, leading to a drop

65    of sequencing reads for low copy DNA/RNA molecules present in a sample [7]. In

66    addition, ligation requires double-stranded DNA (dsDNA) inputs and may result in

67    cross- and self-ligation adaptor byproducts [8].  To overcome these limitations,

68    SMART, a template switching method, was developed and used initially for

69   transcriptome analyses, such as CAGE, RNA-seq (including small RNA-seq),

70   and single-cell RNA-seq [9-12]. By using single-step adapter addition, the

71   SMART technology achieves a much-needed sensitivity to accurately amplify

72   picogram quantities of nucleic acids.

73

74   The SMART method was adapted for preparing NGS libraries from DNA

75   templates in 2014 by tailing an adaptor to the 3' end of a target DNA sequence

76   and later amplifying the sequence by template switching. This modification allows

77   quick preparation of DNA libraries from picogram quantities of DNA molecules [7].

78   Soon, this strategy was applied to ChIP-seq studies with human, mouse and

79   yeast samples [13-19], and it is one of the few currently available protocols for

80   ChIP-seq studies of small cell numbers [20, 21]. Here, a stretch of Ts is added to

81   DNA templates in the tailing step, which is subsequently hybridized to a poly(dA)

82   primer used to copy DNA (**Fig. 1a**). It is conceivable that the poly(dA) primer,

83   however, can lead to signals amplified from non-targeted genomic regions

84   containing consecutive Ts. Indeed, a recent study of SMART ChIP-seq reads

85   revealed a strong bias of base constitution at the 3' end of the sequenced reads

86   that are enriched near long ($\geq$12bp) poly(T/A) containing genomic loci [14]. The

87   authors proposed a computational strategy to reduce this bias by normalizing the

88   ChIP-seq data for the genomic abundance of different polyN tracts, but only

89   achieved partial success [14]. Here, we revisited this problem and demonstrated

90   that the unique features of the falsely amplified reads can be exploited to

91   effectively remove artifact ChIP-seq reads from SMART protocols. We

4

92    implemented this idea in the software SMARTcleaner. Testing multiple published

93    ChIP-seq data, we showed that SMARTcleaner could properly identify and

94    remove artifact reads in both paired-end (PE) and single-end (SE) ChIP-seq data,

95    leading to improved ChIP-seq results.

96

## Results

98    **Strand-specific false priming and amplification at the poly(T/A) sites**

99    When the SMART protocol (or kit) is applied to prepare NGS libraries from DNA

100    fragments, such as those from chromatin immunoprecipitation (IP), there are five

101    steps, 1) 3' T-tailing, 2) annealing of DNA SMART poly(dA) primer to the T-tails,

102    3) primer extension by the SMARTScribe$^{TM}$ reverse transcriptase (RT), 4)

103    template switching and extension by RT using SMART oligo, and 5) PCR-

104    mediated addition of Illumina adapters and subsequent amplification (**Fig. 1a**).

105    As mentioned previously [14], the SMART poly(dA) primers can anneal to poly(T)

106    sequences that are either located within the IP-DNA fragments (**Fig. 1b**) or

107    present in non-target DNA fragments (i.e., the DNA fragments pulled down

108    during IP non-specifically) (**Fig. 1c**). In both cases, the Ts are from genomic

109    sequences and are not added during the T-tailing process. After amplification,

110    sequencing, and read mapping (note that only one strand of the dsDNA is

111    sequenced), ChIP-seq reads from poly(T/A) genomic DNAs, due to false priming

112    and amplification, will accumulate next to the poly(T/A) sites in a clear strand-

113    specific manner because the poly(dA) primers only anneal to the DNA strand

114    containing poly(T). To illustrate this, we examined the reads in a human ChIP-

115    seq sample (Additional file 1: Table S1, Dataset 1, SRR3229031) that was

116    prepared using the Clontech DNA SMART ChIP-seq kit and by PE sequencing

117    [14]. As this particular dataset was obtained from sequencing of control samples

118    (i.e., input DNA), no genomic regions would be expected to show ChIP-seq read

119    enrichment. Indeed, at non-poly(T/A) sites, we did not find accumulations of

120    reads on either "+" or "-" strands (**Fig. 1d**). However, at poly(T/A) sites, we

121    observed that the Read2 of the PE reads were piled up either at the upstream of

122    the poly(T) sites (with respect to the reference "+" strand) (**Fig. 1e**) or at the

123    downstream of the poly(A) sites (**Fig. 1f**), as reported [14]. If SE sequencing had

124    been performed, the accumulation of reads would still be observed, but the

125    precise location information provided by Read2 would not be available (**Fig. 1e,f**),

126    because only Read1 (**Fig. 1a-c)** would be sequenced. Genome-wide analysis of

127    read distribution aggregated over poly(T/A) sites further illustrate these patterns

128    (**Fig. 1g-i**). The width of the peaks indicates the range where the false fragments

129    are located near the poly(T/A) sites (**Fig. 1g-i**).

130

131    **Random false priming and amplification at consecutive and intermittent**

132    **poly(T/A) sites**

133    We reasoned that the SMART poly(dA) primers can anneal to and amplify poly(T)

134    sequences, allowing some degree of mismatch. The PE sequencing data in the

135    SRR3229031 dataset allowed us to identify exactly the ChIP-seq fragments that

136    were artifacts from the poly(T/A) genomic sites, because the Read2 of the

137    fragments would be piled up at the end of poly(T/A) (**Fig. 1e,f;** Additional file 2:

6

138    Figure S1). We should point out that the second reads of the PE sequences

139    submitted to the SRA database have been cut by 10 bp from the 3' end by the

140    authors [14], resulting in a 10 bp gap between the poly(T/A) sites and the end of

141    the Read2 (Additional file 2: Figure S1).

142

143    We counted the numbers of ChIP-seq Read2 that mapped to the 9,698,838

144    poly(A) and 9,796,521 poly(T) sequences containing a minimal of five

145    consecutive As or Ts, respectively, in the human genome (hg38). Like a previous

146    study [14], we found that the median counts for the regions with 5 to 11

147    consecutive A or T were 1, while the median for regions with 12 As or Ts was

148    doubled, indicating that the false priming event occurs primarily at sites with 12 or

149    more consecutive poly(T/A) bases (Additional file 2: Figure S2a; Wilcoxon test, *p*-

150    value < 2.2e-16). Nevertheless, there were large variations at the poly(T/A) sites

151    of the same length, a common phenomenon due to the randomness in primer

152    annealing and sequencing (Additional file 2: Figure S2a). To consider

153    mismatching during priming, we focused on short poly(T/A) sites (≤8bp) that by

154    themselves cannot be efficiently used for false priming but jointly may be. We

155    found that read numbers mapped to two such sequences disrupted by one

156    mismatch nucleotide were significantly reduced, compared to those without

157    disruption, indicating reduced efficiency of false priming (Additional file 2: Figure

158    S1c,d, Figure S2b). Moreover, an insertion of two or three mismatch nucleotides

159    basically abolished false priming (Additional file 2: Figure S2b). In short, our

160    analysis confirmed that false priming occurs significantly at regions containing a

7

161    consecutive sequence of ≥12 As or Ts and the resultant artifact reads should be

162    excluded from ChIP-seq data analysis.

163

164    **SMARTcleaner: identification and cleaning of falsely primed fragments**

165    Based on the above information of the false priming event in SMART ChIP-seq

166    studies (**Fig. 1**, Additional file 2: Figure S2), we developed a computational tool,

167    SMARTcleaner, to remove the ChIP-seq artifact signals. It has two modes (PE

168    mode and SE mode) to accommodate the two sequencing options during ChIP-

169    seq. In PE mode, a genome (FASTA) sequence file and ChIP-seq read

170    alignment files (in bam format) are taken as input, and "cleaned" bam files are

171    generated with the reads predicted from false priming removed and saved in the

172    "noise" bam files. In SE mode, it takes a list of consecutive and interrupted

173    poly(T/A) genomic sites (Additional file 2: Figure S2), and bam files, and outputs

174    cleaned bam files and noise bam files. The software is publicly available through

175    github (https://github.com/dzhaobio/SMARTcleaner).

176

177    In PE mode, our tool removes ChIP-seq read pairs whose second reads mapped

178    to poly(T/A) (see Methods). Analysis of pileup reads at individual poly(A/T) sites

179    (**Fig. 2a,b**) and total read counts across all poly(A/T) sites (**Fig. 2c,d**)

180    demonstrated clearly that reads from false priming in the SRR3229031 dataset

181    were effectively identified and successfully removed by SMARTcleaner.

182    Furthermore, applying the SMARTcleaner to ChIP-seq data from libraries

183    constructed using a ligation method [14], we found that < 0.002% of PE reads

8

184   were mistakenly removed, indicating that the PE mode is highly accurate. By

185   comparison, artifact reads in the SMART-based data could be successfully

186   removed, while their percentages (11-20%) varied among the different DNA

187   shearing methods used for fragmentation (**Fig. 2e**). In addition, for the SMART-

188   based data, the ChIP-seq fragment sizes calculated from the noise bam files

189   were 21-43 bp shorter on average than those in the clean bam files, as expected,

190   since the genomic poly(T/A) sequences were within ChIP fragments while tailed

191   Ts were added to the ends of ChIP fragments. This observation is consistent with

192   previous finding [14].

193

194   In SE mode, the SMARTcleaner identifies and removes artifact reads by

195   comparing read distributions in the "+" and "-" strands near individual poly(T/A)

196   sites, because false priming leads to reads accumulated in only one of the two

197   strands (**Fig. 1**).  To demonstrate its performance, we treated the above PE

198   ChIP-seq reads as SE reads, by analyzing the Read1 data only.  Again, analysis

199   of pileup reads at individual poly(T/A) sites (**Fig. 3a,b**) and read counts

200   aggregated over genome wide poly(T/A) sites (**Fig. 3c,d**) demonstrated that most

201   artifact reads were removed effectively. However, the SE mode appeared less

202   robust than the PE mode, because it mistakenly removed ~0.8% of reads in the

203   ligation-based ChIP-seq data (**Fig. 3e**). The percentages of reads that were

204   removed by the SE mode for the SMART-based datasets were similar to those

205   using the PE mode (**Fig. 3e**).

206

9

207 In terms of computational efficiency, we tested both PE and SE modes on a PC

208 (Intel(R) Xeon(R) CPU E5-2609 0 @ 2.40GHz, 32Gb memory, CentOS Linux

209 release 7.3.1611). It took 30 min to clean 94 million reads in PE mode and 16

210 min to clean 47 million reads in SE mode, benchmarking with the SRR3229031

211 dataset. The PE mode requires more memory than the SE mode because the

212 former reads the entire genome sequence into memory (for fast query) and

213 keeps track of the end coordinates of Read2 at the genomic poly(T/A) sites.

214

215 **Evaluation of SMARTcleaner with published histone modification ChIP-seq**

216 **datasets**

217 To demonstrate the value of our tool and importance of removing artifact reads

218 from false priming in the analysis of SMART ChIP-seq data, we first applied the

219 SMARTcleaner to a public ChIP-seq dataset (Additional file 1: Table S1, Dataset

220 2) that studied H3K4me3 histone modification in HeLa cells using seven methods

221 for preparing sequencing libraries from low-input IP DNAs, including SMART

222 method [13]. The study also generated a PCR-free dataset as a gold standard

223 reference, including three replicates using 100 ng DNA as starting material. For

224 the other seven protocols, the starting material was either 1 ng or 0.1 ng, each

225 with five replicates [13]. The original study was designed for comparing the

226 performance of different ChIP-seq library preparation methods, but this dataset is

227 ideal for evaluating our tool for three reasons. First, its gold standard data can be

228 used for clearly evaluating artifacts introduced in PCR amplification. Second, the

229 dataset is valuable for evaluating the effect of initial DNA inputs on false priming

10

230     and amplification. Third, the known enrichment of H3K4me3 peaks at promoter

231     regions [22] can be used as a metric to measure the impact of falsely called

232     peaks.

233

234     In our test below, as a benchmark we chose the data from PCR-free method and

235     Ascel2S method, which were consistently ranked at the top by multiple criteria in

236     the original study [13]. Since the ChIP-seq libraries were sequenced by the

237     single-end method, we applied SE mode to the alignment files, including control

238     samples. Similar to the above finding in **Fig. 3e**, only a small percentage of ChIP-

239     seq reads were removed by SMARTcleaner from the ligation-based datasets, 0.3%

240     on average. For SMART-derived dataset, the average percentage was 3.0% for

241     1 ng and 5.3% for 0.1 ng starting DNA material (Additional file 2: Figure S3a).

242     Next, we randomly sampled 6 millions of reads for each sample for calling

243     H3K4me3 peaks using the software MACS2 [23], by the same criteria. We found

244     that before read cleaning 12.1% and 17.1% of the H3K4me3 peaks, called from

245     the 1 ng and 0.1 ng SMART protocols respectively, overlapped with poly(T/A)

246     sites, but after cleaning the overlaps dropped to 6.2% and 8.1%, comparable to

247     the numbers for PCR-free and Ascel2S samples (Additional file 2: Figure S3b).

248     This result indicates that not all peaks in poly(T/A) sites are artifacts. The greater

249     percentages of removed reads and peak overlaps with poly(T/A) sites for the 0.1

250     ng than the 1 ng dataset are consistent with the assumption of increased false

251     priming when the input DNA material is lower, due to a reduced number of

252     genuine target DNA templates. In addition, the percentages of H3K4me3 peaks

11

253    mapping to promoters increased by 3.7% (1 ng) and 4.1% (0.1 ng) after cleaning

254    reads in the SMART derived datasets, while the change (0.14%) is negligible for

255    the PCR free and Ascel2S samples (Additional file 2: Figure S3c).

256

257    We also compared the SMART ChIP-seq peaks to the H3K4me3 peaks from

258    PCR-free samples, using the peaks (n= 20,262) present in all three PCR-free

259    datasets as the reference. The mean sensitivity (i.e., % PCR-free peaks detected

260    in SMART) was 89.68% and 89.61% in pre- and post-cleaning samples (1ng

261    DNA), indicating no difference in sensitivity. Same was observed for the samples

262    using 0.1ng starting DNA material (Additional file 2: Figure S3d). However, the

263    specificity (% SMART peaks found in PCR-free peaks) was increased from 89.25%

264    to 90.42% for samples with 1ng DNA and from 87.11% to 89.85% for samples

265    with 0.1ng DNA after cleaning the noise (Additional file 2: Figure S3e), indicating

266    that the cleaning process improved the peak quality.

267

268    Next, we directly compared the pre- and post-cleaning H3K4me3 peak lists. The

269    total number of peaks dropped for both SMART samples after cleaning (**Fig. 4a**),

270    but the change for 0.1 ng SMART sample was significant larger than that for 1 ng

271    one (**Fig. 4b**), clearly suggesting that with lower amounts of input DNA, more

272    false peaks would be called from the artifact reads (**Fig. 4c**). In support of this,

273    we observed that the 0.1 ng pre-cleaning SMART samples had the largest

274    percentages (on average 64.3%) of peaks located near the poly(T/A) sites (**Fig.**

275    **4d**). When compared to the peaks called for the PCR-free data, 51.9% (0.1 ng)

276    and 35.1% (1 ng) of the peaks unique to the pre-cleaning SMART samples

277    overlapped, significantly smaller than the percentages for peaks either shared

278    with or unique to post-cleaned data (**Fig. 4e**). Similarly, the percentages of

279    H3K4me3 peaks (44.4% and 39.8%) located to promoters for the peaks unique

280    to pre-cleaning samples were significantly lower than the numbers for the other

281    two groups of peaks (**Fig. 4f**). As an orthogonal measurement, we analyzed

282    transcription factor (TF) motifs in the H3K4me3 peak regions. The TATA box and

283    CAAT box, two well-known general promoter TF motifs [24], and the ETS motif

284    [25], were  the most enriched motifs in the H3K4me3 peaks. In all cases, their

285    occurrences in the peaks detected only in the pre-cleaning samples were

286    significantly lower (**Fig. 4g-i**). In contrast, the RLR1 motif, which basically

287    consists of poly(T), was only enriched in the peaks unique to the pre-cleaning

288    samples (**Fig. 4j**). Finally, we examined the ChIP-seq read densities and

289    aggregated read profiles for the three groups of H3K4me3 peaks, unique to pre-

290    or post-cleaning samples, or shared (**Fig. 4k**). The peaks unique to the post-

291    cleaning samples had about 2x stronger (both 1 ng and 0.1 ng samples) ChIP-

292    seq signals in the PCR-free and Ascel2S data than the peaks unique to the pre-

293    cleaning samples, indicating that the latter peaks were very likely derived from

294    PCR amplification and thus enriched for artifacts (**Fig. 4k**). Taken together, these

295    results indicate that the reads removed by SMARTcleaner are true artifacts and

296    its application can improve the quality of peaks identified from ChIP-seq analysis,

297    resulting in better biological findings.

298

**Evaluation of SMARTcleaner with published transcription factor ChIP-seq datasets**

We were especially interested in how the inclusion of artifact reads may affect peaks identified from TF ChIP-seq studies. Therefore, we reanalyzed a previously published Olig2 ChIP-seq dataset (Additional file 1: Table S1, Dataset 3) and compared our results to the original publication [18]. We found that 16% of the original peaks (3,251 of 20,283) overlapped with the poly(T/A) sites, with some peaks exhibiting typical features of false amplification (**Fig. 5a**). We also noticed that the authors applied a combination of very stringent criteria to filter peaks, perhaps in an effort to limit peaks from false priming. Thus, we tried less stringent criteria to obtain a new set of peaks (n=25,179) from the pre-cleaning alignment files and included it in our comparison (see Methods). Next, we used the SMARTcleaner SE mode to clean the alignment files and obtained a list of post-cleaning peaks (n=23,289). A comparison of the three lists of peaks is shown in **Fig. 5b**, from which we defined four groups of peaks (Additional file 2: Figure S4): "TP", or true positive, called by all methods; "FP", or false positive, called by the original study and present in the pre-cleaning sample only; "FN", or false negative, removed by the original study only; and "TN", or true negative, removed in the original study and by SMARTcleaner. Intersections of the four groups of peaks with poly(T/A) sites showed that 92.9% of TN peaks and 94.3% of FP peaks overlapped with poly(T/A) sites, compared to 12.7% of TP peaks and 5.3% of FN peaks (**Fig. 5c**), indicating that the original study not only included some artifact peaks but also filtered out some true peaks. This was

14

322    supported by a comparative analysis of the ChIP-seq read intensities, with reads

323    from false priming present in both the ChIP sample and input control (FP and TN

324    in **Fig. 5d,e**). This analysis also showed that the FN group represented true

325    peaks filtered out by the authors by using overly strict criteria (**Fig. 5d,e**).

326

327    To further test the cleaning effect, we included a Olig2 ChIP-seq dataset that was

328    independently generated from neural stem cells using a non-SMART protocol

329    [26]. We found that 86.2% and 91.8% of the pre-cleaning and post-cleaning

330    peaks were detected by the non-SMART method, respectively. Moreover, among

331    the four groups of peaks, 93.8% and 83% of TP and FN peaks were present in

332    the non-SMART peaks, respectively, in contrast to 8.7% and 6.2% for the TN and

333    FP groups, respectively, indicating that false peaks were removed by our clearing

334    process. This result was supported by the patterns in the read density heatmaps

335    and profiles (**Fig. 5d,e**).

336

337    In addition, motif analysis demonstrated that the top four motifs enriched in the

338    TP and FN peaks were the same TF motifs (Atoch1, NF1, Tcf12 and Olig2)

339    reported in the original study [18]. However, the top motifs for the TN and FP

340    groups were RLR1, TA repeat, GAGA repeat, CTCF and Myf5, which seem

341    irrelevant to Olig2 function (**Fig. 5f**).

342

343    In short, our analysis of the Olig2 ChIP-seq data further supports the value of our

344    newly developed SMARTcleaner tool, and illustrates the need for appropriately

15

345    removing noise and artifacts from false priming in TF ChIP-seq studies that use

346    the SMART protocol.

347

348    **Prevalence of artifact reads from false priming and amplification in SMART-**

349    **based ChIP-seq datasets**

350    To determine if false priming and amplification is a common problem in SMART-

351    based ChIP-seq libraries, we collected and analyzed all such datasets except a

352    clinical one that is not publicly accessible [15] (Additional file 1: Table S1; see

353    Methods). These ChIP-seq data were carried out in human [13-16], mouse [17,

354    18], and yeast samples [19]. All but two of the datasets were analyzed by single-

355    end sequencing [14, 15]. Our analysis showed that all available datasets

356    contained an average of 8.5% (2.7% ~19.6%) reads that were likely derived from

357    false priming, regardless of the amount of input DNA (from 0.1 ng to 10 ng DNA)

358    or cell numbers (from 10 to 100 millions) (Additional file 1: Table S1).

359

360    **Discussion**

361    The SMART ChIP-seq kit uses the template switching method to improve the

362    efficiency of library construction, which is especially suitable for analyzing

363    samples with very low amounts of input DNA [7]. Consistent with a recent report

364    [14], we show that the protocol, however, can introduce significant noise to ChIP-

365    seq data, due to the annealing of DNA SMART poly(dA) primers to non-targeted

366    genomic regions containing ≥ 12 Ts or As. The artifact reads have distinct

367    features (**Fig. 1**, Additional file 2: Figure S2) that are exploited by the

16

368     SMARTcleaner tool developed in this study. Using multiple published ChIP-seq

369     datasets, we demonstrated convincingly that our tool can successfully remove

370     the artifact reads arising from false priming and amplification of the SMART

371     poly(dA) primers. It works for both PE and SE ChIP-seq reads (**Fig. 2**, **Fig. 3**),

372     and outputs both cleaned alignment files and noise, which can be loaded into a

373     genome browser for inspecting the cleaning effects visually. SMARTcleaner also

374     provides some running options and helper tools to prepare the files required for

375     the cleaning process. Currently SMARTcleaner does not deal with biases

376     introduced by other factors, such as DNA shearing method etc. [5], but users can

377     easily adapt this tool to their ChIP-seq analytic pipelines and develop it further.

378

379     We have examined all currently available public datasets that were obtained

380     using the DNA SMART ChIP-seq kit, and found that the false priming issue is

381     prevalent, regardless of the amount of input DNA material or cell numbers

382     (Additional file 1: Table S1). While the artifact cannot be easily removed by data

383     normalization, strict filtering in peak calling, or a simple exclusion of peaks

384     located at poly(A/T) sites, our study suggests that the false priming issue

385     becomes less severe when a large amount of DNA is used as the starting

386     material for ChIP library preparation. Conceivably, the concern can also be

387     alleviated if high affinity antibodies are used to significantly enrich target DNA

388     templates in the input material. Based on our survey of all available datasets, we

389     have the following recommendations to users of the SMART ChIP-seq kit to

390     exploit its full potential. First, one should use a sufficient amount of DNA as the

17

391     starting templates, whenever possible. Second, the T-tailing step in the SMART

392     ChIP-seq protocol should be optimized. Third, sequence the NGS libraries using

393     the PE method and clean the ChIP-seq reads using the PE mode of

394     SMARTcleaner. Forth, if the libraries have already been sequenced using the SE

395     method, clean the ChIP-seq reads using the SE mode of SMARTcleaner.

396     Alternatively, one can consider to use other ChIP-seq library preparation

397     methods that can also handle low-input DNA [13, 20, 21].

398

399     **Conclusions**

400     False priming and amplification occur at poly(T/A) genomic sites due to the use

401     of poly(dA) primers in SMART-based ChIP-seq library construction. Reads from

402     subsequent false amplification and sequencing are strand-specific and can be

403     effectively removed by our SMARTcleaner tool, leading to improvement in peak

404     calling, and downstream data analysis and interpretation.

405

406     **Methods**

407     **ChIP-seq datasets and read processing**

408     The SMART ChIP-seq kit is a promising but relatively new protocol for analyzing

409     small amount of chromatin materials. We searched for ChIP-seq datasets that

410     used this kit in the GEO and by Google and found one publication in 2015 [18],

411     two in 2016 [13, 19], and four in 2017 [14-17]. Among the seven publications, six

412     have made their data publicly accessible (Additional file 1: Table S1). The

413    seventh is a clinical study and the corresponding data have not been released,

414    possibly due to protection of privacy [15]. In the alignment of ChIP-seq reads

415    derived from the SMART protocols, the first three bases were trimmed from the

416    first read (Read1). In all datasets, replicates were analyzed independently. To

417    facilitate comparison with the original studies, we used the same versions of

418    software as in the original publication when applicable.

419    *Dataset 1*

420    The first dataset is actually a ChIP-seq of input DNAs from HCT116 cells and

421    HeLa-S3 because the DNA templates were not enriched with any antibodies. It

422    contained seven sets of paired-end sequencing data, which we downloaded from

423    the NCBI SRA database (SRP071830) [14]. Three libraries were constructed

424    using the DNA SMART ChIP-Seq kit (Clontech, #634865), with the others by

425    "standard" ligation-based method. Reads were mapped to the human genome

426    (hg38) using Bowtie2 (v2.2.3) [27], using default parameters with the maximum

427    fragment length for valid paired-end reads set to 2000. Only uniquely mapped

428    reads were kept for further analyses, after duplicate reads were removed using

429    the Picard tool -- MarkDuplicates (v2.3.0,

430    http://broadinstitute.github.io/picard/index.html). To mimic single-end sequencing,

431    we generated SE bam files by extracting the first reads from the PE bam files

432    (samtools view -h -f 64).

433    *Dataset 2*

434    The H3K4me3 ChIP experiments were done with 56 million HeLa cells in 56

435    ChIP reactions [13]. The ChIP DNA was combined into a single pool and then

19

436    divided into seven aliquots for different library preparation methods and the PCR-

437    free method. Libraries starting from either 1 ng or 0.1 ng ChIP DNA were

438    generated. Reads were aligned to the hg38 human reference genome using

439    Bowtie (v1.2.1) [28]. Only uniquely mapped reads were used for analysis, with

440    duplicate reads removed by samtools (v0.1.19) [29]. To call peaks, we randomly

441    subsampled 6 million mapped reads for each sample, as done in the original

442    study [13] and used the MACS2 (v2.1.0) [23] with $q$ value < 0.05. Motif analysis

443    was done using the HOMER (v4.7) [30].

444    **_Dataset 3_**

445    The Olig2 ChIP-seq was carried out with 10 million neural stem cells (NSCs)

446    derived from embryonic (E14.5) CD-1 mice. The libraries were constructed using

447    the DNA SMART ChIP-seq kit and sequenced by the single-end method on an

448    Illumina HiSeq2000 sequencer [18]. The dataset was downloaded from the GEO

449    database (GEO: GSE74646). Reads were aligned to the mouse reference

450    genome (mm10) using bowtie (v1.2.1). Only uniquely mapped reads were used

451    for analysis, with duplicate reads removed using samtools (v0.1.19). Peaks were

452    called using the MACS (v1.4.2) and filtered by $p$ value < $10^{-5}$, fold enrichment > 5,

453    and tag number > 15. When the filter was set to the same as used in the original

454    paper ($p$ value < $10^{-9}$, fold enrichment > 5, and tag number > 20), we obtained

455    essentially the same peaks that were called in the original study. Peak motif

456    analysis was done using HOMER (v4.7) [30].

20

457 *Dataset 4*

458 The H3K4me1 ChIP-seq was obtained with 10 million SUM159 cells. H3K4me1

459 ChIP-seq libraries were constructed using the DNA SMART ChIP Seq Kit

460 (Clontech) with 10ng ChIP DNA (NCBI GEO: GSE87424) [16]. Raw fastq

461 sequences were downloaded from the GEO and processed with the same

462 methods as the original study.

463 *Dataset 5*

464 The ChIP-seq experiments of H3K27ac histone modification and c-MYC were

465 performed with FACS-sorted Eph4 cells. Libraries were constructed using the

466 Clontech DNA Smart Chipseq kit (Clontech, #634866), and pooled for

467 sequencing (NCBI GEO: GSE98004) [17]. Raw fastq sequences were

468 downloaded from the GEO and processed as the original study.

469 *Dataset 6*

470 The last dataset was from a yeast study [19]. DNA–RNA immunoprecipitation

471 and deep sequencing (DRIP-seq) was done with S9.6 monoclonal antibody in

472 100 million yeast cells. We downloaded the alignment files from European

473 Nucleotide Archive (ENA) website (PRJEB8021) and yeast reference genome

474 from the UCSC genome browser [31].

475

476 **SMARTcleaner**

477 The SMARTcleaner tool was developed in Perl under the MIT license after

478 analysis of the characteristics of ChIP-seq reads derived from false priming and

21

479    amplification. Two modes, PE mode and SE mode, were implemented based on

480    the sequencing methods used in ChIP-seq data.

481    ***PE mode***

482    When sequenced in PE method, the second reads of the falsely primed

483    fragments will pile up upstream of the poly(T) sites or the downstream of the

484    poly(A) sites (**Figure 1e,f**), allowing two mismatch insertions (Additional file 2:

485    Figure S2). SMARTcleaner will go through a sorted (by coordinates) alignment

486    file and find read pairs with the second read at the left end of poly(T) sites or at

487    the right end of poly(A) sites (Additional file 2: Figure S5). It will keep tracking the

488    number of such fragments at each position of a poly(T/A) site. When this number

489    is over a threshold (default: 1) predefined for false amplification, all read pairs

490    ending in the same position will be considered as artifacts and placed to the new

491    alignment file ("noise bam file"). In the meantime, the original bam file subtracting

492    the artifact reads will be saved as a cleaned bam file.

493    ***SE mode***

494    When ChIP-seq is sequenced in SE method, the false reads will be clustered

495    upstream of poly(T) sites or downstream of poly(A) sites of the reference genome

496    (**Fig. 1**), up to two mismatches (Additional file 2: Figure S2). SMARTcleaner first

497    examines the reads in the flanking regions (by default 2kb) of all poly(T/A) sites

498    to decide the size of the region containing falsely amplified fragments. For reads

499    on "+" strand, the distance is calculated from the left ends of reads to the left

500    ends of poly(T) sites (Additional file 2: Figure S6a) or the right ends of poly(A)

501    sites (Additional file 2: Figure S6b). For reads on "-" strand, the distance is

22

502    calculated from the right ends of reads to the left ends of poly(T) sites (Additional

503    file 2: Figure S6a) or the right ends of poly(A) sites (Additional file 2: Figure S6b).

504    Based on the distribution of the distances, SMARTcleaner automatically

505    determines the window size at poly(T/A) sites for sampling, or a user can

506    manually set it according to the read distribution at the poly(T/A) sites (**Fig. 1h,i**).

507    A bed file containing the resampling regions will be generated. Next, it will go

508    through the reads at each of those regions, check if the potentially artifact reads

509    outnumber (default 2x) those in the unaffected opposite strand, and finally

510    resample the artifact reads, if necessary, according to the read numbers in the

511    opposite strand (Additional file 2: Figure S7a,b). For the genomic regions with

512    overlapping poly(T) and poly(A) sites, the tool will process the poly(T/A) sites

513    based on the order of their appearance in the reference genome (Additional file 2:

514    Figure S7c).

515

516    For SE mode, a list of poly(T/A) sites is needed. We included a helper command

517    to identify such regions in a genome. To estimate the range for resampling reads,

518    we implemented another helper command in our tool for this purpose. Users can

519    also directly set a range for resampling based on their knowledge of their

520    datasets or the fragment distribution around the poly(T/A) sites.

521

522    **List of abbreviations**

523    ChIP-seq: chromatin immunoprecipitation and deep sequencing

524    NGS : next-generation sequencing

525 PE: paired-end

526 SE: single-end

527 SMART: switching mechanism at the 5' end of the RNA transcript

528

## Declarations

529

530 **Ethics approval and consent to participate**

531 Not applicable

532 **Consent for publication**

533 Not applicable

534 **Availability of data and material**

535 The datasets supporting the conclusions of this article are available in the NCBI

536 SRA: SRP071830 (Dataset 1), NCBI SRA: SRP067250 (Dataset 2), NCBI GEO:

537 GSE74646 (Dataset 3), NCBI GEO: GSE87424 (Dataset 4), NCBI GEO:

538 GSE98004 (Dataset 5), and European Nucleotide Archive (ENA): PRJEB8021

539 (Dataset 6). SMARTcleaner is publicly available under MIT license at github

540 (https://github.com/dzhaobio/SMARTcleaner).

541 **Competing interests**

542 The authors declare that they have no competing interests.

543 **Funding**

**Authors' contributions**

D. Zhao developed the tool, performed the analyses, and drafted the manuscript.

D. Zheng contributed ideas, wrote the manuscript, and supervised the study. All

authors read and approved the final manuscript.

**Acknowledgements**

**References**

1.      Kahvejian A, Quackenbush J, Thompson JF. What would you do if you
        could sequence everything? Nat Biotech 2008;26:1125-1133.

2.      Shendure J, Aiden EL. The expanding scope of DNA sequencing. Nat
        Biotech 2012;30:1084-1094.

3.      Rabbani B, Nakaoka H, Akhondzadeh S, Tekin M, Mahdieh N. Next
        generation sequencing: implications in personalized medicine and
        pharmacogenomics. Molecular BioSystems 2016;12:1818-1830.

25

567   4.   Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of

568        next-generation sequencing technologies. Nat Rev Genet 2016;17:333-

569        351.

570   5.   van Dijk EL, Jaszczyszyn Y, Thermes C. Library preparation methods for

571        next-generation sequencing: Tone down the bias. Experimental Cell

572        Research 2014;322:12-20.

573   6.   Head SR, Komori HK, LaMere SA, Whisenant T, Van Nieuwerburgh F,

574        Salomon DR, Ordoukhanian P. Library construction for next-generation

575        sequencing: overviews and challenges. Biotechniques 2014;56:61-64, 66,

576        68, passim.

577   7.   Turchinovich A, Surowy H, Serva A, Zapatka M, Lichter P, Burwinkel B.

578        Capture and Amplification by Tailing and Switching (CATS). RNA Biology

579        2014;11:817-828.

580   8.   Raabe CA, Tang TH, Brosius J, Rozhdestvensky TS. Biases in small RNA

581        deep sequencing data. Nucleic Acids Res 2014;42:1414-1426.

582   9.   Tang DTP, Plessy C, Salimullah M, Suzuki AM, Calligaris R, Gustincich S,

583        Carninci P. Suppression of artifacts and barcode bias in high-throughput

584        transcriptome analyses utilizing template switching. Nucleic Acids

585        Research 2013;41:e44-e44.

586   10.  Ramskold D, Luo S, Wang Y-C, Li R, Deng Q, Faridani OR, Daniels GA,

587        Khrebtukova I, Loring JF, Laurent LC *et al*. Full-length mRNA-Seq from

588        single-cell levels of RNA and individual circulating tumor cells. Nat Biotech

589        2012;30:777-782.

590   11.   Picelli S, Bjorklund AK, Faridani OR, Sagasser S, Winberg G, Sandberg

591         R. Smart-seq2 for sensitive full-length transcriptome profiling in single

592         cells. Nat Meth 2013;10:1096-1098.

593   12.   Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M,

594         Chen Y, Zhao X, Schmidl C, Suzuki T *et al*. An atlas of active enhancers

595         across human cell types and tissues. Nature 2014;507:455-461.

596   13.   Sundaram AYM, Hughes T, Biondi S, Bolduc N, Bowman SK, Camilli A,

597         Chew YC, Couture C, Farmer A, Jerome JP *et al*. A comparative study of

598         ChIP-seq sequencing library preparation methods. BMC Genomics

599         2016;17:816.

600   14.   Vardi O, Shamir I, Javasky E, Goren A, Simon I. Biases in the SMART-

601         DNA library preparation method associated with genomic poly dA/dT

602         sequences. PLOS ONE 2017;12:e0172769.

603   15.   Vong JSL, Tsang JCH, Jiang P, Lee W-S, Leung TY, Chan KCA, Chiu

604         RWK, Lo YMD. Single-Stranded DNA Library Preparation Preferentially

605         Enriches Short Maternal DNA in Maternal Plasma. Clinical Chemistry

606         2017;63:1031-1037.

607   16.   Zawistowski JS, Bevill SM, Goulet DR, Stuhlmiller TJ, Beltran AS,

608         Olivares-Quintero JF, Singh D, Sciaky N, Parker JS, Rashid NU *et al*.

609         Enhancer Remodeling during Adaptive Bypass to MEK Inhibition Is

610         Attenuated by Pharmacologic Targeting of the P-TEFb Complex. Cancer

611         Discovery 2017;7:302-321.

27

612    17.    Frey WD, Chaudhry A, Slepicka PF, Ouellette AM, Kirberger SE,

613            Pomerantz WCK, Hannon GJ, dos Santos CO. BPTF Maintains Chromatin

614            Accessibility and the Self-Renewal Capacity of Mammary Gland Stem

615            Cells. Stem Cell Reports 2017;9:23-31.

616    18.    Dong X, Chen K, Cuevas-Diaz Duran R, You Y, Sloan SA, Zhang Y, Zong

617            S, Cao Q, Barres BA, Wu JQ. Comprehensive Identification of Long Non-

618            coding RNAs in Purified Cell Types from the Brain Reveals Functional

619            LncRNA in OPC Fate Determination. PLOS Genetics 2015;11:e1005669.

620    19.    Wang IX, Grunseich C, Chung YG, Kwak H, Ramrattan G, Zhu Z, Cheung

621            VG. RNA–DNA sequence differences in Saccharomyces cerevisiae.

622            Genome Research 2016;26:1544-1554.

623    20.    Dahl JA, Gilfillan GD. How low can you go? Pushing the limits of low-input

624            ChIP-seq. Briefings in Functional Genomics 2017:elx037-elx037.

625    21.    Skene PJ, Henikoff S. An efficient targeted nuclease strategy for high-

626            resolution mapping of DNA binding sites. eLife 2017;6:e21856.

627    22.    Liu X, Wang C, Liu W, Li J, Li C, Kou X, Chen J, Zhao Y, Gao H, Wang H

628            *et al*. Distinct features of H3K4me3 and H3K27me3 chromatin domains in

629            pre-implantation embryos. Nature 2016;537:558-562.

630    23.    Zhang Y, Liu T, Meyer C, Eeckhoute J, Johnson D, Bernstein B, Nusbaum

631            C, Myers R, Brown M, Li W *et al*. Model-based Analysis of ChIP-Seq

632            (MACS). Genome Biology 2008;9:R137.

633 24. Cooper SJ, Trinklein ND, Anton ED, Nguyen L, Myers RM.

634    Comprehensive analysis of transcriptional promoter structure and function

635    in 1% of the human genome. Genome Research 2006;16:1-10.

636 25. Sharrocks AD. The ETS-domain transcription factor family. Nat Rev Mol

637    Cell Biol 2001;2:827-837.

638 26. Mateo JL, van den Berg DLC, Haeussler M, Drechsel D, Gaber ZB, Castro

639    DS, Robson P, Lu QR, Crawford GE, Flicek P *et al*. Characterization of

640    the neural stem cell gene regulatory network identifies OLIG2 as a

641    multifunctional regulator of self-renewal. Genome Research 2015;25:41-

642    56.

643 27. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2.

644    Nat Meth 2012;9:357-359.

645 28. Langmead B, Trapnell C, Pop M, Salzberg S. Ultrafast and memory-

646    efficient alignment of short DNA sequences to the human genome.

647    Genome Biology 2009;10:R25.

648 29. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G,

649    Abecasis G, Durbin R, Subgroup GPDP. The Sequence Alignment/Map

650    format and SAMtools. Bioinformatics 2009;25:2078-2079.

651 30. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX,

652    Murre C, Singh H, Glass CK. Simple combinations of lineage-determining

653    transcription factors prime cis-regulatory elements required for

654    macrophage and B cell identities. Mol Cell 2010;38:576-589.

655    31.    Casper J, Zweig AS, Villarreal C, Tyner C, Speir ML, Rosenbloom KR,

656            Raney BJ, Lee CM, Lee BT, Karolchik D *et al*. The UCSC Genome

657            Browser database: 2018 update. Nucleic Acids Research 2017:gkx1020-

658            gkx1020.

659    32.    Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz

660            G, Mesirov JP. Integrative genomics viewer. Nat Biotech 2011;29:24-26.

661    33.    Ye T, Krebs AR, Choukrallah M-A, Keime C, Plewniak F, Davidson I, Tora

662            L. seqMINER: an integrated ChIP-seq data interpretation platform. Nucleic

663            Acids Research 2011;39:e35.

664

665


666    **Figure legends**

667    **Fig. 1. Strand-specific amplification of non-targeted sequences at poly(T/A)**

668    **sites in the SMART ChIP-seq analysis**. **a.** Flowchart of the SMART ChIP-seq

669    procedure at non-poly(T/A) sites, adapted from the user manual of the kit

670    (http://www.clontech.com/xxclt_ibcGetAttachment.jsp?cItemId=99449). **b,c**.

671    Modified flowcharts to show annealing of the SMART poly(dA) primers to non-

672    tailed Ts within targeted (**b**) or non-targeted (**c**) DNA templates, leading to

673    strand-specific amplification at poly(T) sites. For poly(A) sites, false amplification

674    occurs to the opposite strand. **d-f**. ChIP-seq read densities at three randomly

675    picked non-poly(T/A) and poly(T/A) sites. The data is from SRR3229031

676    (Additional file 1: Table S1, Dataset 1), and Integrative Genomics Viewer (IGV)

677    [32] is used to show the ChIP-seq reads from paired-end (PE) or single-end (SE)

678    sequencing. For PE, read1 and read2 are shown as pairs, with reads mapped to

679    "+" and "-" strands in red and blue, respectively. For SE, only Read1 (extracted

680    from PE data) is shown. **g-i**. Aggregated read distribution at non-poly(T/A) and

681    poly(T/A) sites. In h and i, poly(T/A) sites were defined as those with ≥ 12

682    consecutive T or A in the human reference genome. To define non-poly(T/A)

683    sites, we first selected genomic regions that are > 4 kb in length and > 1kb away

684    from poly(T/A) sites, and then take the 2kb regions around the middle points. In

685    total, we got 301,474 non-poly(T/A) sites, 338,568 poly(T) sites, and 336,703

686    poly(A) sites. Refer to the Method section (SE mode, Additional file 2: Figure S6)

687    for the calculation of read distribution.

688

689    **Fig. 2. SMARTcleaner in PE mode**. **a**. PE reads mapped to a poly(T) and a

690    poly(A) locus before (raw) and after cleaning. **b**. A genomic region showing the

691    read densities before and after cleaning. The "called peaks" refer to pre-cleaning

692    peaks called using MACS2. **c,d**. Genome-wide read distribution at poly(T/A) sites

693    before (red and blue lines) and after (green lines) cleaning. **e**. Percentages of

694    removed reads at poly(T/A) sites in each sample. The samples from left to right

695    are SRR3229030, SRR3286889, SRR3286890, SRR3286891, SRR3229031,

696    SRR3286910, and SRR3286911 (Additional file 1: Table S1, Dataset 1).

697

698    **Fig. 3. SMARTcleaner in SE mode**. **a**. Two examples showing the cleaning

699    results of SE mode at one poly(T) and one poly(A) locus. **b**. Cleaning result in a

700    genomic region. **c,d**. Genome-wide reads distribution near the poly(T/A) sites

31

701  before (red and blue lines) and after (green lines) cleaning. **e**. Percentages of

702  removed reads at poly(T/A) sites in samples prepared by ligation or SMART

703  protocols. The sample order is the same as in Fig. 2e.

704

705  **Fig. 4**. **Evaluate SMARTcleaner with H3K4me3 ChIP-seq data**. **a**. Numbers of

706  pre- and post-cleaning H3K4me3 peaks. **b**. Change of peak numbers after

707  cleaning. **c**. Numbers of peaks shared or unique to pre-cleaning ("uniqPre") or

708  post-cleaning ("uniqPost") data **d**. Overlap of peaks with poly(T/A) sites. **e**.

709  Overlap of peaks with gold standard peaks. **f**. Peaks at the promoter regions (2kb

710  around TSS). **g-j**. Percentages of peaks with each of the four enriched TF motifs.

711  **k**. Read densities and average profiles for peaks shared by or unique to pre- and

712  post-cleaning data. Reads counts were extracted using seqMINER [33] from 6

713  million reads randomly sampled from individual samples. Heatmaps were drawn

714  using R package pheatmap, with peaks as row and sorted by read densities. In

715  **a-j**, each point represents a replicate sample.

716

717  **Fig. 5**. **Evaluate SMARTcleaner with a TF ChIP-seq data**. **a**. An example of

718  false peaks in the original list of Olig2 ChIP-seq peaks. The track of "called peak"

719  shows peaks provided by the authors. **b**. Venn diagram showing the peak

720  overlaps from three methods: the original peaks from the authors, the peaks

721  called before cleaning, and the peaks called after cleaning. When counting the

722  overlapping peaks, we could get two different numbers depending on which set

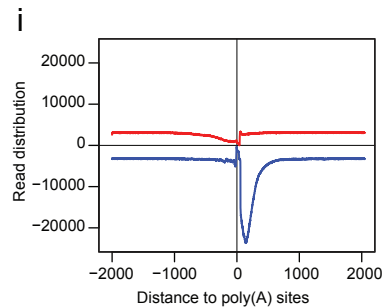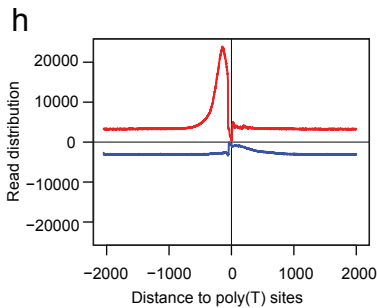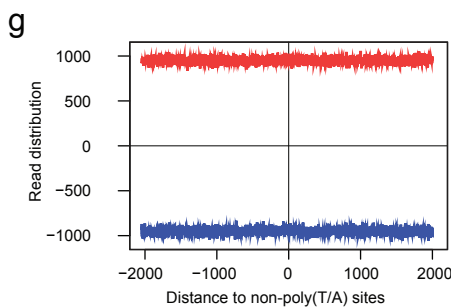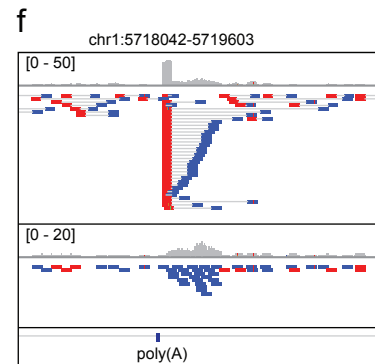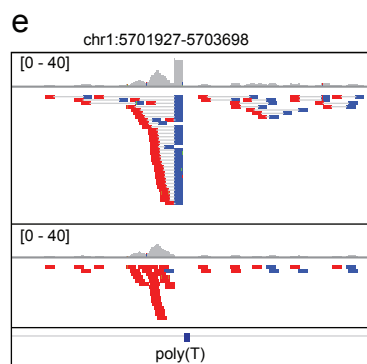723  of peaks is used to report the number (one peak in one set may overlap more
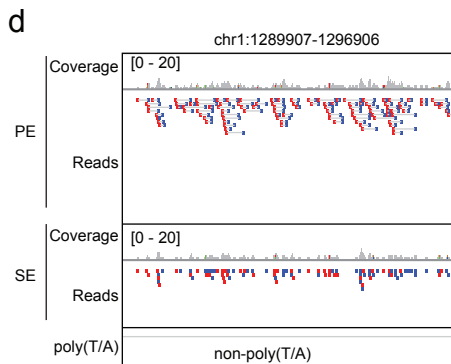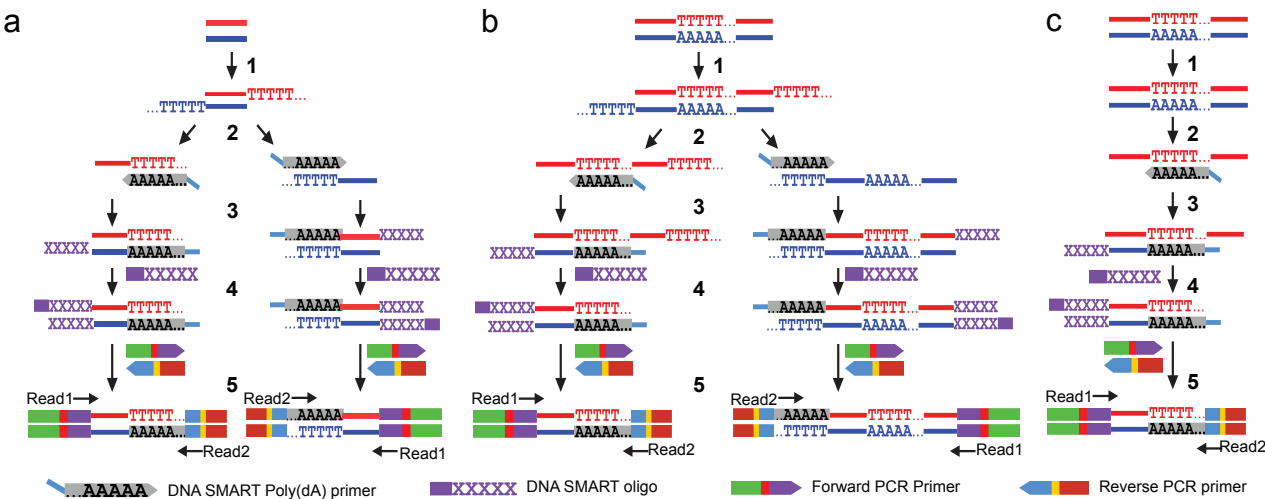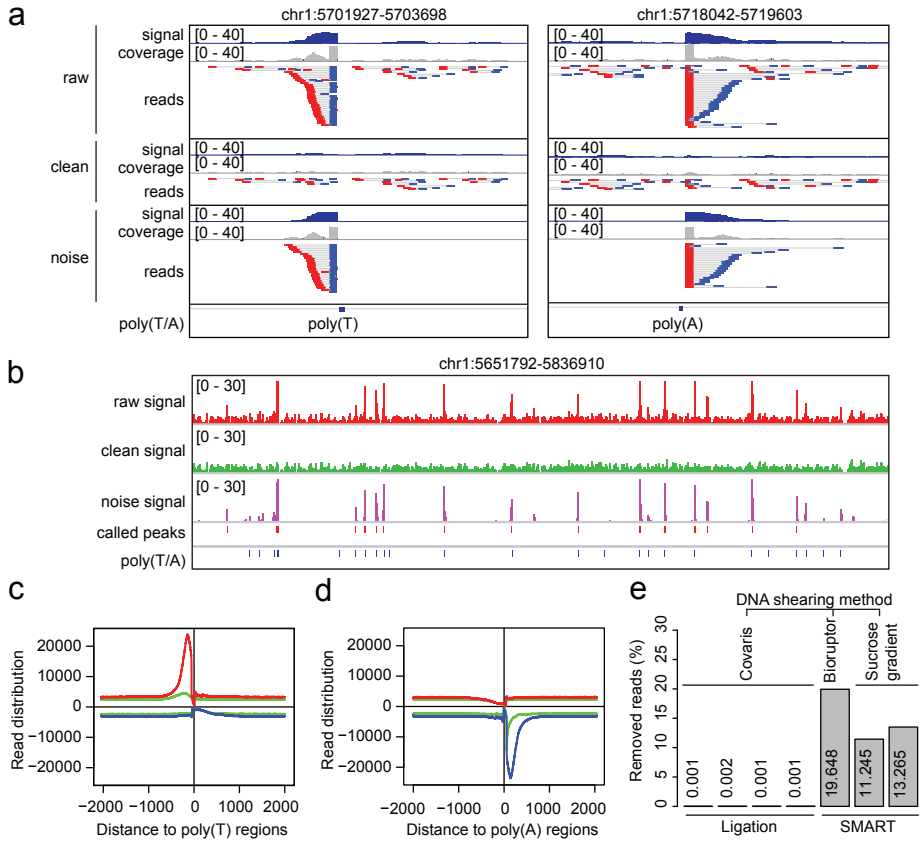
724    than one peak in another set). We reported the smaller number here. **c**. Peaks

725    overlapping with poly(T/A) sites. **d,e**. Read densities and average counts at the

726    four selected groups of peaks, computed by sampling 5 million reads. An Olig2

727    ChIP-seq data (right) from non-SMART method was also analyzed. **f**. Top

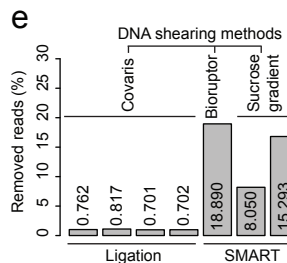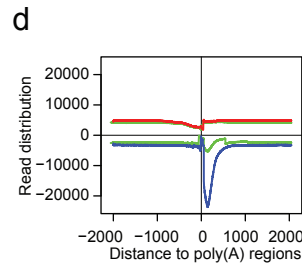728    enriched motifs by HOMER [30].

729

730    **Description of additional data files**

731    Additional file 1: Supplementary Table S1.

732    Additional file 2: Supplementary Figure S1–S7.

**a**

**b**

**c**

DNA SMART Poly(dA) primer

DNA SMART oligo

Forward PCR Primer

Reverse PCR primer

**d** chr1:1289907-1296906

**e** chr1:5701927-5703698

**f** chr1:5718042-5719603

Coverage [0 - 20]

Coverage [0 - 40]

Coverage [0 - 50]

PE

Reads

Coverage [0 - 20]

Coverage [0 - 40]

Coverage [0 - 20]

SE

Reads

poly(T/A) non-poly(T/A)

poly(T)

poly(A)

**g**

**h**

**i**

Read distribution

Distance to non-poly(T/A) sites

Distance to poly(T) sites

Distance to poly(A) sites

**a**

chr1:5701927-5703698     chr1:5718042-5719603

raw — signal [0 - 40], coverage [0 - 40], reads

clean — signal [0 - 40], coverage [0 - 40], reads

noise — signal [0 - 40], coverage [0 - 40], reads

poly(T/A) — poly(T)    poly(A)

**b**

chr1:5651792-5836910

raw signal [0 - 30]

clean signal [0 - 30]

noise signal [0 - 30]

called peaks

poly(T/A)

**c** Read distribution vs Distance to poly(T) regions

**d** Read distribution vs Distance to poly(A) regions

**e** Removed reads (%) — DNA shearing method

Covaris: 0.001, 0.002, 0.001, 0.001 (Ligation)
Bioruptor: 19.648; Sucrose gradient: 11.245, 13.265 (SMART)

**a**

chr1:5701927-5703698 | chr1:5718042-5719603

raw
- signal [0 - 20]
- coverage [0 - 20]
- reads (+)
- reads (−)

clean
- signal [0 - 20]
- coverage [0 - 20]
- reads (+)
- reads (−)

noise
- signal [0 - 20]
- coverage [0 - 20]
- reads (+)
- reads (−)

poly(T/A) — poly(T) | poly(A)

**b**

chr1:5651792-5836910

- raw signal [0 - 20]
- clean signal [0 - 20]
- noise signal [0 - 20]
- called peaks
- poly(T/A)

**c**

Read distribution vs Distance to poly(T) regions

**d**

Read distribution vs Distance to poly(A) regions

**e**

Removed reads (%) — DNA shearing methods

Covaris: 0.762, 0.817, 0.701, 0.702
Bioruptor: 18.890
Sucrose gradient: 8.050, 15.293

Ligation | SMART

**a**

chr1:13040343-13044223

Olig2 ChIP
- signal [0 - 21]
- reads(+)
- reads(-)

IgG input
- signal [0 - 10]
- reads(+)
- reads(-)

poly(T/A)

called peak

**b** peaks in the original paper

- 62
- FP 1068
- 3
- TP 19211
- TN 917
- FN 3984
- 26
- pre-cleaning
- post-cleaning

**c** % of peaks overlapping poly(T/A) sites

TP  FN  TN  FP

**d**

SMART — Olig2 ChIP (pre-cleaning, post-cleaning), IgG input (pre-cleaning, post-cleaning)

Non-SMART — Olig2 ChIP

TP, FN, TN, FP

**e** Average read coverage

TP, FN, TN, FP

-2kb  0  +2kb

**f**

| Group | Motif | TF | P-value |
|---|---|---|---|
| TP | | Atoh1 | 1e-904 |
| | | NF1 | 1e-850 |
| | | Tcf12 | 1e-802 |
| | | Olig2 | 1e-518 |
| FN | | Atoh1 | 1e-180 |
| | | Tcf12 | 1e-152 |
| | | NF1 | 1e-143 |
| | | Olig2 | 1e-136 |

| Group | Motif | TF | P-value |
|---|---|---|---|
| TN | | RLR1 | 1e-231 |
| | | TA repeat | 1e-79 |
| | | GAGA repeat | 1e-34 |
| | | CTCF | 1e-25 |
| FP | | RLR1 | 1e-494 |
| | | TA repeat | 1e-186 |
| | | GAGA repeat | 1e-57 |
| | | Myf5 | 1e-28 |