1    **STARRPeaker: Uniform processing and accurate identification of STARR-seq active**

2    **regions**

3

4    Donghoon Lee[1], Manman Shi[2], Jennifer Moran[2], Martha Wall[2], Jing Zhang[1,3], Jason Liu[3,]

5    Dominic Fitzgerald[2], Yasuhiro Kyono[2], Lijia Ma[2,4], Kevin P White[2,5]*, Mark Gerstein[1,3,6,7]*

6

7    1 Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT

8    06520, USA

9    2 Institute for Genomics and System Biology, University of Chicago, Chicago, IL, 60637, USA

10    3 Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT

11    06520, USA

12    4 School of Life Sciences, Westlake University, Hangzhou, 310024, China

13    5 Tempus Labs, Inc. Chicago IL 60654, USA

14    6 Department of Computer Science, Yale University, New Haven, CT 06520, USA

15    7 Department of Statistics and Data Science, Yale University, New Haven, CT 06520, USA

16

17    * Corresponding author

18    E-mail: pi@gersteinlab.org

19

20

21 **Abstract**

22 High-throughput reporter assays, such as self-transcribing active regulatory region sequencing

23 (STARR-seq), allow for unbiased and quantitative assessment of enhancers at a genome-wide

24 level. Recent advances in STARR-seq technology have employed progressively more complex

25 genomic libraries and increased sequencing depths, to assay larger sized regions, up to the entire

26 human genome. These advances necessitate a reliable processing pipeline and peak-calling

27 algorithm. Most STARR-seq studies have relied on chromatin immunoprecipitation sequencing

28 (ChIP-seq) processing pipeline to identify peaks. However, there are key differences in STARR-

29 seq versus ChIP-seq data: STARR-seq uses transcribed RNA to measure enhancer activity,

30 making determining the basal transcription rate important. Furthermore, STARR-seq coverage is

31 non-uniform, overdispersed, and often confounded by sequencing biases such as GC content and

32 mappability. Moreover, here, we observed a clear correlation between RNA thermodynamic

33 stability and STARR-seq readout, suggesting that STARR-seq might be sensitive to RNA

34 secondary structure and stability. Considering these findings, we developed STARRPeaker: a

35 negative binomial regression framework for uniformly processing STARR-seq data. We applied

36 STARRPeaker to two whole human genome STARR-seq experiments; HepG2 and K562. Our

37 method identifies highly reproducible and epigenetically active enhancers across replicates.

38 Moreover, STARRPeaker outperforms other peak callers in terms of identifying known

39 enhancers. Thus, our framework optimized for processing STARR-seq data accurately

40 characterizes cell-type-specific enhancers, while addressing potential confounders.

41

42 **Keywords**: STARR-seq, peak caller, enhancer, non-coding

43

## **Introduction**

The transcription of eukaryotic genes is precisely coordinated by an interplay between cis-regulatory elements. For example, enhancers and promoters serve as platforms for transcription factors (TF) to bind and interact with each other, and their interactions are often required to initiate transcription[1,2]. Enhancers, which are often located distantly from the transcribed gene body itself, play critical roles in the upregulation of gene transcription. Enhancers are cell-type specific and can be epigenetically activated or silenced to modulate transcriptional dynamics over the course of development. Enhancers can be found upstream or downstream of genes, or even within introns[3–5]. They function independent from their orientation, do not necessarily regulate the closest genes, and sometimes regulate multiple genes at once[6,7]. In addition, several recent studies have demonstrated that some promoters – termed E-promoters – may act as enhancers of distal genes[8,9].

Unlike protein-coding genes, enhancers do not yet have a well-characterized consensus sequence. Therefore, identifying enhancers in an unbiased fashion is challenging. The non-coding territory occupies over 98% of the genome landscape, making the search space very broad. Moreover, the activity of enhancers depends on the physiological condition and epigenetic landscape of the cellular environment, complicating the fair assessment of enhancer function.

Previously, putative regulatory elements were computationally predicted, indirectly, by profiling DNA accessibility (using DNase-seq, FAIRE-seq, and ATAC-seq) as well as histone modifications (ChIP-seq) that are linked to regulatory functions[10–12]. More recently, researchers have developed high-throughput episomal (exogenous) reporter assays to directly measure

3

67    enhancer activity across the whole genome, specifically massively parallel reporter assays

68    (MPRA)[13,14] and self-transcribing active regulatory region sequencing (STARR-seq)[15,16]. These

69    assays allow for quantitative assessment of enhancer activity in a high-throughput fashion.

70

71    In STARR-seq, candidate DNA fragments are cloned downstream of a reporter gene into the 3′

72    untranslated region (UTR). After transfecting the plasmid pool into host cells, one can measure

73    the regulatory potential by high-throughput sequencing of the 3′ UTR of the expressed reporter

74    gene mRNA. These exogenous reporters enable accurate and unbiased assessment of enhancer

75    activity at the whole genome level, independent of chromatin context. Unlike MPRA – which

76    utilizes barcodes – STARR-seq produces self-transcribed RNA fragments that can be directly

77    mapped onto the genome. The activities of enhancers are measured by comparing the amount of

78    RNA produced from the input DNA library. STARR-seq has several technical advantages over

79    MPRA. Library construction is relatively simple because barcodes are not needed. In addition,

80    candidate enhancers are cloned instead of synthesized, allowing the assay to test extended

81    sequence contexts (>500 bp) for enhancer activity, which studies have shown to be critical for

82    functional activity[17]. Importantly, STARR-seq can be scaled to the whole genome level for

83    unbiased scanning for functional elements. However, scaling STARR-seq to the human genome

84    is still very challenging, primarily due to its massive size. A more complex genomic DNA

85    library, a higher sequencing depth, and increased transfection efficiency are required to cover the

86    whole human genome[16], which could ultimately introduce biases.

87

88    Processing of STARR-seq is somewhat similar to chromatin immunoprecipitation sequencing

89    (ChIP-seq), where protein-crosslinked DNA is immunoprecipitated and sequenced. A typical

90   ChIP-seq processing pipeline identifies genomic regions over-represented by sequencing tags in

91   a ChIP sample compared to a control sample. STARR-seq data is compatible with most ChIP-

92   seq peak callers. Hence, previous studies on STARR-seq have largely relied on peak calling

93   software developed for ChIP-seq such as MACS2[16,18,19]. However, one must be cautious using

94   ChIP-seq peak callers, at least without re-tuning default parameters optimized for processing

95   transcription factor ChIP-seq[20].

96

97   In this paper, we describe key differences in the processing of STARR-seq versus ChIP-seq data.

98   Due to increased complexity of the genomic screening library and sequencing depth

99   requirements, STARR-seq coverage is highly non-uniform. This leads to a lower signal-to-noise

100   ratio than a typical ChIP-seq experiment and makes estimating the background model more

101   challenging, which could ultimately lead to false positives peaks. In addition, STARR-seq

102   measures more of a continuous activity similar to quantification in RNA-seq than a discrete

103   binding event. Therefore, STARR-seq peaks should be further evaluated using a notion of

104   activity score. These differences necessitate a unique approach to processing STARR-seq data.

105

106   We propose an algorithm optimized for processing and identifying functionally active enhancers

107   from STARR-seq data, which we call STARRPeaker. This approach statistically models the

108   basal level of transcription, accounting for potential confounding factors, and accurately

109   identifies reproducible enhancers. We applied our method to two whole human STARR-seq

110   datasets and evaluated its performance against previous methods. We also compared an R

111   package, BasicSTARRseq, developed to process peaks from the first STARR-seq data[15], which

112   models enrichment using a binomial distribution. We benchmarked our peak calls against known

113    human enhancers. Thus, our findings support that STARRPeaker will be a useful tool for

114    uniformly processing STARR-seq data.

115

116    **Materials and Methods**

117

118    *Precise measurement of STARR-seq coverage*

119    We binned the genome using a sliding window of length, *l*, and step size, *s*. Based on the average

120    size of the STARR-seq library, we defined a 500 bp window length with a 100 bp step size to be

121    the default parameter. Based on generated genomic bins, we calculated the coverage of both

122    STARR-seq input and output mapped to each bin. For calculating the sequence coverage, other

123    peak callers and many visualization tools commonly use the start position of the read[15,21,22].

124    However, given that the average sizes of the fragments inserted in STARR-seq libraries were

125    approximately 500 bp, we expected that the read coverage using the start position of read may

126    shift the estimate of the summit of signal and dilute the enrichment. Some peak callers have used

127    read densities of forward and reverse strand separately to overcome this issue[23,24]. To precisely

128    measure the coverage of STARR-seq input and output, we first inferred the size of the fragment

129    insert from paired-end reads and used the center of the fragment insert, instead of start position

130    of the read, to calculate coverage. For inferring the size of fragment insert, we first strictly

131    filtered out reads that were not properly paired and chimeric. Chimeric alignments are reads that

132    cannot be linearly aligned to a reference genome, implying a potential discrepancy between the

133    sequenced genome and the reference genome and indicative of structural variation[25]. We also

134    filtered out read pairs that had a fragment insert size less than $l_{max}$ and greater than $l_{min}$. By

135    default, we filtered out fragment insert sizes less than 100 bp and greater than 1,000 bp. After

6

136     filtering out spurious read-pairs, we estimated the center of the fragment insert and counted the

137     fragment depth for each genomic bin. We compared the coverage calculated using the start of

138     read against the center of fragment insert and observed both a shift in the location of enrichment

139     summit and a difference in enrichment level (**Figure 1**).

140

141     ***Controlling for potential systemic bias in sequencing and STARR-seq library preparation***

142     STARR-seq measures the ratio of transcribed RNA to DNA for a given test region and

143     determines whether the test region can facilitate transcription at a higher rate than the basal level.

144     This is based on the assumption that the basal transcriptional level stays relatively constant

145     across the genome and the transcriptional rate is a reflection of the regulatory activity of a test

146     region. However, this may not always be true, and one needs to consider potential systemic

147     biases when analyzing the result. Unlike ChIP-seq where both the experiment and input controls

148     are from the same DNA origin, STARR-seq experiments measure the regulatory potential from

149     the abundance of transcribed RNA, which adds a layer of complexity. For example, RNA

150     structure and co-transcriptional folding might potentially influence the readout of STARR-seq

151     experiments[26]. Single-stranded RNA starts to fold upon transcription and the resulting RNA

152     structure might influence the measurement of regulatory activity. Previously, researchers

153     suggested a potential linkage between RNA secondary structure and transcriptional regulation[27].

154     In addition, the resulting transcribed RNA undergoes a series of post-transcriptional regulation,

155     and RNA stability might play a critical role. Moreover, previous reports have shown that the

156     degradation rates vary significantly across the genome and RNA degradation rates are the main

157     determinant of cellular RNA levels[28]. Furthermore, RNA stability correlates with

158     functionality[29,30].

159

160    There are also intrinsic sequencing biases in library preparation. A genome-wide reporter library

161    is made from randomly sheared genomic DNA, but DNA fragmentation is often non-random[31].

162    Studies have also suggested that epigenetic mechanisms and CpG methylation may influence

163    fragmentation[32]. Furthermore, the isolated polyadenylated RNAs are reverse transcribed and

164    PCR is amplified before sequenced, and this process can further confound the sequenced

165    candidate fragments.

166

167    To unbiasedly test for the regulatory activity, a model needs to control for these potential

168    systemic biases inherent to generating STARR-seq data. As we expected, we observed that

169    STARR-seq coverage for both input and output are confounded by potential sequencing bias

170    (**Figure 2**). Notably, STARR-seq coverage significantly correlated with GC content (PCC 0.61;

171    P-val 1E-299), mappability (PCC 0.45; P-val 2.9E-148), and RNA thermodynamic stability

172    (PCC -0.55; P-val 0). Hence, to unbiasedly identify the activity peaks from STARR-seq, we

173    developed a model that accounts for variability of tested candidate fragments.

174

175    ***Accurate modelling of STARR-seq coverage using negative binomial regression***

176    To model the fragment coverage data from STARR-seq using discrete probability distribution,

177    we assumed that each genomic bin is independent and identically distributed, as specified in

178    Bernoulli trials[33]. That is, each test fragment can only map to a single fixed-length bin. Therefore,

179    we only considered a non-overlapping subset of bins for modeling and fitting the distribution.

180    We also excluded bins not covered by any genomic input or normalized input coverage was less

181    than a minimum quantile $t_{min}$, since these regions do not have sufficient power to detect

8

182    enrichment. We simulated and fitted various discrete probability distributions to STARR-seq

183    coverage. We observed that the STARR-seq coverage data was overdispersed and fitted the best

184    with negative binomial distribution (**Figure 3A**). We also noticed a slight negative enrichment,

185    indicating that some candidate fragments can silence the basal transcriptional activity. A Q-Q

186    plot of simulated coverage further demonstrated that the negative binomial model provides the

187    best fit for the data (**Figure 3B**).

188

189    *Peak caller*

190    To accurately model the ratio of STARR-seq sequence coverage (RNA) to input sequence

191    coverage (DNA) while controlling for potential confounding factors, we applied a negative

192    binomial regression. The overview of our model is outlined in **Figure 4**. Our model starts by

193    fitting an analytical distribution to the observed fragment coverage across each genomic bin. In

194    doing so, we use covariates to model expected counts in the form of multiple regression. Once

195    regression coefficients are estimated from a set of data, we can evaluate the likelihood of

196    observing the fragment count for each bin and assign p-values. Ultimately, bins with significant

197    enrichments are selected based on an adjusted p-values threshold, and they are fine-tuned to the

198    summit of the peak fragment enrichment.

199

200    Let Y be a vector of STARR-seq output (RNA) coverage, then $y_i$ for $1 \leq i \leq n$ denotes the

201    number of RNA fragments from STARR-seq experiment mapped to the $i$-th bin from the total of

202    $n$ genomic bins. Let $t_i$ be the number of input library (DNA) mapped to the $i$-th bin. We define

203    $X$ be the matrix of covariates where $\vec{x_i}$ is the vector of covariates corresponding to the $i$-th bin,

204    and $x_{ij}$ is the $j$-th covariate for the $i$-th bin.

9

205

206 *Negative binomial distribution*

207 A negative binomial distribution, which arises from a Gamma-Poisson mixture, can be

208 parametrized as follows[34–36] (see Supplementary Methods for derivation).

209

$$f_Y(y_i|\mu_i,\theta) = \frac{\Gamma(y_i+\theta)}{\Gamma(y_i+1)\cdot\Gamma(\theta)} \cdot \left(\frac{\theta}{\theta+\mu_i}\right)^{\theta} \cdot \left(\frac{\mu_i}{\theta+\mu_i}\right)^{y_i}$$

210

211 A negative binomial is a generalization of a Poisson regression that allows the variance to be

212 different from the mean, shaped by the dispersion parameter $\theta$. The variance for the NB2 model

213 is given as

214

$$\sigma^2 = \mu + \frac{\mu^2}{\theta}$$

215

216 We assume that the majority of genomic bins will have a basal level of transcription, and the

217 count of RNA fragments at each $i$-th bin follows the traditional negative binomial (NB2)

218 distribution. The expected fragment counts, $E(y_i)$, represents the mean incidence, $\mu_i$.

219

$$y_i \sim NB(\mu_i,\theta)$$

$$E(y_i) = \mu_i$$

220

221 *Negative binomial regression model*

10

222    The regression term for the expected RNA fragment count can be expressed in terms of a linear

223    combination of explanatory variables, a set of $m$ covariates ($\vec{x}$). We use the input library variable

224    $t_i$ as one covariate. For simplicity, we denote $t_i$ as $x_{0i}$ hereafter.

225

$$\ln \mu_i = \beta_0 x_{0i} + \beta_1 x_{1i} + \cdots + \beta_m x_{mi}$$

$$\mu_i = \exp(\beta_0 x_{0i} + \beta_1 x_{1i} + \cdots + \beta_m x_{mi})$$

$$\mu_i = \exp(\vec{x_i}^\intercal \beta)$$

226

227    Alternatively, instead of using the input library variable $t_i$ as one covariate, we can directly use it

228    as an offset variable. One advantage of using the input variable as an "exposure" to the RNA

229    output coverage is that it allows us to directly model the basal transcription rate (the ratio of

230    RNA to DNA) as a rate response variable. More details on this alternative parametrization are

231    described in the Supplementary Methods.

232

233    *Maximum-likelihood estimation*

234    We fit the model and estimate regression coefficients using the maximum likelihood method,

235    where log-likelihood function is shown as follows.

236

$$\mathcal{L}_{NB}(\mu|y,\theta) = \sum_{i=1}^{n} y_i \ln\left(\frac{\mu_i}{\theta + \mu_i}\right) + \theta \ln\left(\frac{\theta}{\theta + \mu_i}\right) + \ln\left(\frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)}\right)$$

237

238    Substituting $\mu_i$ with the regression term, the log-likelihood function can be parametrized in terms

239    of regression coefficients, $\beta$.

11

240

$$\mathcal{L}_{NB}(\beta|y,\theta) = \sum_{i=1}^{n} y_i \ln\left(\frac{e^{\overrightarrow{x_i}^{\top}\beta}}{\theta + e^{\overrightarrow{x_i}^{\top}\beta}}\right) + \theta \ln\left(\frac{\theta}{\theta + e^{\overrightarrow{x_i}^{\top}\beta}}\right) + \ln\left(\frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1)\cdot\Gamma(\theta)}\right)$$

241

242    We can determine the maximum likelihood estimates of the model parameters by setting the first

243    derivative of the log-likelihood with respect to $\beta$, the gradient, to zero, and there is no analytical

244    solution for $\hat{\beta}$. Numerically, we iteratively solve for the regression coefficients $\beta$ and the

245    dispersion parameter $\theta$, alternatively, until both parameters converge.

246

247    *Estimation of P-value*

248    Finally, we calculate a P-value based on the fitted value of the $i$-th bin from the cumulative

249    distribution function of negative binomial distribution, and we assign false discovery rate using

250    Benjamini & Hochberg method[37].

251

$$P\text{-}value = \Pr(x \geq y_i) = 1 - CDF(x = y_i - 1)$$

$$= 1 - \sum_{i=0}^{\hat{y}_i - 1} \binom{\hat{y}_i + \theta - 1}{\hat{y}_i} \frac{\theta}{\theta + \hat{y}_i}^{\hat{y}_i} (1 - \frac{\theta}{\theta + \hat{y}_i})^{\theta}$$

252

253    **Source code and data availability**

254    We implemented the method described in this article as a Python software package called

255    STARRPeaker. The software package can be downloaded, installed, and readily used to call

256    peaks from any STARR-seq dataset. The STARRPeaker package, as well as source code and

257    documentation, is freely available at: http://github.com/gersteinlab/starrpeaker. Data used in the

258    analysis will be made available from the Gene Expression Omnibus for public use.

259    DNase-seq and ChIP-seq data used for the analysis is publicly available from the ENCODE

260    portal (https://www.encodeproject.org/). The specific accession codes used for the analysis are

261    listed in Supplementary Table S3. GC content was downloaded from the UCSC Genome

262    Browser (http://hgdownload.cse.ucsc.edu/gbdb/hg38/bbi/gc5BaseBw/), and the mappability

263    track was created using gem-library software[38] with a k-mer size of 100 bp and the reference

264    human genome build hg38.

265

266    **<u>Results</u>**

267    We applied our peak calling algorithm to two whole human genome STARR-seq experiments,

268    K562 and HepG2, utilizing origin of replication-based (ORI) plasmids. Using this dataset, we

269    evaluated the quality and characteristics of identified enhancers as well as the performance of the

270    peak caller by comparing to external enhancer datasets.

271

272    ***Accurate identification of highly reproducible enhancers***

273    To evaluate the quality of enhancers identified from STARRPeaker, we uniformly called peaks

274    from the whole human genome STARR-seq dataset using methods previously used to identify

275    enhancers from STARR-seq data, namely BasicSTARRseq and MACS2, using recommended

276    settings. We first compared the level of epigenetic profile enrichment around the peaks. We

277    observed higher enrichment of DNase hypersensitive sites, as well as more distinct double-peak

278    patterns of H3K27ac and H3K4me1, using STARR-seq versus BasicSTARRseq or MACS2

279    (**Figure 5**). We also aggregated the transcription factor binding sites assayed by ChIP-seq around

13

280    peaks, and we observed significant enrichment of transcription factor binding events compared

281    to peaks identified by other methods. Furthermore, we compared STARRPeaker peaks and

282    others to previously characterized enhancers by CAGE[39], MPRA[17,40], and STARR-seq[19] in

283    HepG2 or K562 cell line (**Figure 6**). We observed a higher fraction of STARRPeaker peaks

284    overlap with external datasets.

285

286    **<u>Discussion</u>**

287    We developed a statistically rigorous analysis pipeline for STARR-seq data in a software

288    package named STARRPeaker. STARRPeaker has several key improvements over previous

289    peak identification methods including (1) accurate quantification of STARR-seq coverage based

290    on inferred fragment size from paired-end reads; (2) use of a negative binomial distribution to

291    account for overdispersion in bin counts; and (3) modeling of STARR-seq coverage as a function

292    of input and potential confounding variables in STARR-seq signal. We applied our method to

293    two whole human genome ORI-STARR-seq datasets and demonstrated that it can unbiasedly

294    identify a set of STARR-seq-positive regions better than previous methods. The STARR-seq

295    peaks were enriched with epigenetic marks relevant to enhancers and overlapped better with

296    previously known enhancers than previous methods.

297

298    To completely understand how noncoding regulatory elements can modulate transcriptional

299    programs in human, STARR-seq active regions must be further characterized and validated

300    within the cellular context. Currently, CRISPR-based screens are limited to a small number of

301    selected targets. Our method can aid in prioritize candidate regions in unbiased fashion to

302    maximize the functional characterization efforts.

14

303

**Funding**

306

**Acknowledgements**

312

**Author Contributions**

314    D.L., M.S., K.W., and M.G. conceived the project. D.L. and M.G. drafted the manuscript. D.L.

315    developed the STARRPeaker software package. M.S., J.M., M.W., D.F., Y.K., and L.M.

316    performed experimental work. M.W. performed experimental validation. D.L., J.Z., and J.L.

317    performed the downstream analysis. M.G. and K.W. provided funding and supervised the project.

318

319

320 **References**

321  1.    Muerdter, F., Boryń, Ł. M. & Arnold, C. D. STARR-seq — Principles and applications.

322      *Genomics* **106**, 145–150 (2015).

323  2.    Yáñez-Cuna, J. O., Kvon, E. Z. & Stark, A. Deciphering the transcriptional cis-regulatory

324      code. *Trends Genet.* **29**, 11–22 (2013).

325  3.    Lettice, L. A. *et al.* A long-range Shh enhancer regulates expression in the developing

326      limb and fin and is associated with preaxial polydactyly. *Hum. Mol. Genet.* **12**, 1725–1735

327      (2003).

328  4.    Banerji, J., Rusconi, S. & Schaffner, W. Expression of a beta-globin gene is enhanced by

329      remote SV40 DNA sequences. *Cell* **27**, 299–308 (1981).

330  5.    Sagai, T., Hosoya, M., Mizushina, Y., Tamura, M. & Shiroishi, T. Elimination of a long-

331      range cis-regulatory module causes complete loss of limb-specific Shh expression and

332      truncation of the mouse limb. *Development* **132**, 797–803 (2005).

333  6.    Melo, C. A. *et al.* eRNAs Are Required for p53-Dependent Enhancer Activity and Gene

334      Transcription. *Mol. Cell* **49**, 524–535 (2013).

335  7.    Sanyal, A., Lajoie, B. R., Jain, G. & Dekker, J. The long-range interaction landscape of

336      gene promoters. *Nature* **489**, 109–13 (2012).

337  8.    Dao, L. T. M. *et al.* Genome-wide characterization of mammalian promoters with distal

338      enhancer functions. *Nat. Genet.* **49**, 1073–1081 (2017).

339  9.    Diao, Y. *et al.* A tiling-deletion-based genetic screen for cis-regulatory element

340      identification in mammalian cells. *Nat. Methods* **14**, 629–635 (2017).

341  10.   Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and

342      characterization. *Nat. Methods* **9**, 215–216 (2012).

343   11.   Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure

344         through genomic segmentation. *Nat. Methods* **9**, 473–476 (2012).

345   12.   Sethi, A. *et al.* A cross-organism framework for supervised enhancer prediction with

346         epigenetic pattern recognition and targeted validation. *bioRxiv* 385237 (2018).

347         doi:10.1101/385237

348   13.   Patwardhan, R. P. *et al.* Massively parallel functional dissection of mammalian enhancers

349         in vivo. *Nat. Biotechnol.* **30**, 265–270 (2012).

350   14.   Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in

351         human cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277

352         (2012).

353   15.   Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by

354         STARR-seq. *Science* **339**, 1074–7 (2013).

355   16.   Liu, Y. *et al.* Functional assessment of human enhancer activities using whole-genome

356         STARR-sequencing. *Genome Biol.* **18**, 219 (2017).

357   17.   Klein, J. C. *et al.* A systematic evaluation of the design, orientation, and sequence context

358         dependencies of massively parallel reporter assays. *bioRxiv* 576405 (2019).

359         doi:10.1101/576405

360   18.   Johnson, G. D. *et al.* Human genome-wide measurement of drug-responsive regulatory

361         activity. *Nat. Commun.* **9**, 5317 (2018).

362   19.   Rathert, P. *et al.* Transcriptional plasticity promotes primary and acquired resistance to

363         BET inhibition. *Nature* **525**, 543–547 (2015).

364   20.   Koohy, H., Down, T. A., Spivakov, M. & Hubbard, T. A comparison of peak callers used

365         for DNase-Seq data. *PLoS One* **9**, e96303 (2014).

366  21.  Uren, P. J. *et al.* Site identification in high-throughput RNA-protein interaction data.

367       *Bioinformatics* **28**, 3013–20 (2012).

368  22.  Strbenac, D., Armstrong, N. J. & Yang, J. Y. H. Detection and classification of peaks in 5'

369       cap RNA sequencing data. *BMC Genomics* **14 Suppl 5**, S9 (2013).

370  23.  Zhang, Y. *et al.* Model-based Analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137

371       (2008).

372  24.  Kharchenko, P. V, Tolstorukov, M. Y. & Park, P. J. Design and analysis of ChIP-seq

373       experiments for DNA-binding proteins. *Nat. Biotechnol.* **26**, 1351–1359 (2008).

374  25.  Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,

375       2078–2079 (2009).

376  26.  Lai, D., Proctor, J. R. & Meyer, I. M. On the importance of cotranscriptional RNA

377       structure formation. *RNA* **19**, 1461–1473 (2013).

378  27.  Ringnér, M. & Krogh, M. Folding free energies of 5'-UTRs impact post-transcriptional

379       regulation on a genomic scale in yeast. *PLoS Comput. Biol.* **1**, e72 (2005).

380  28.  Rabani, M. *et al.* Metabolic labeling of RNA uncovers principles of RNA production and

381       degradation dynamics in mammalian cells. *Nat. Biotechnol.* **29**, 436–42 (2011).

382  29.  Yang, E. *et al.* Decay rates of human mRNAs: correlation with functional characteristics

383       and sequence attributes. *Genome Res.* **13**, 1863–72 (2003).

384  30.  Tani, H. *et al.* Genome-wide determination of RNA stability reveals hundreds of short-

385       lived noncoding transcripts in mammals. *Genome Res.* **22**, 947–56 (2012).

386  31.  Poptsova, M. S. *et al.* Non-random DNA fragmentation in next-generation sequencing. *Sci.*

387       *Rep.* **4**, 4532 (2014).

388  32.  Lazarovici, A. *et al.* Probing DNA shape and methylation state on a genomic scale with

18

389        DNase I. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 6376–81 (2013).

390   33.   Papoulis, A. & Athanasios. Probability, random variables and stochastic processes. *New*

391        *York McGraw-Hill, 1984, 2nd ed.* (1984).

392   34.   Hilbe, J. M. *Negative Binomial Regression.* (Cambridge University Press, 2011).

393        doi:10.1017/CBO9780511973420

394   35.   Cameron, A. C. A. & Trivedi, P. K. *Regression Analysis of Count Data.* (Cambridge

395        University Press, 2013). doi:10.1017/CBO9781139013567

396   36.   Hilbe, J. M. *Modeling Count Data.* (Cambridge University Press, 2014).

397        doi:10.1017/CBO9781139236065

398   37.   Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and

399        Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B* **57**, 289–300 (1995).

400   38.   Derrien, T. *et al.* Fast Computation and Applications of Genome Mappability. *PLoS One* **7**,

401        e30377 (2012).

402   39.   Kawaji, H., Kasukawa, T., Forrest, A. & Carninci, P. The FANTOM 5 collection, a data

403        series underpinning mammalian transcriptome atlases in diverse cell types. *Sci. Data*

404        2016–2018 (2017). doi:10.1038/sdata.2017.113

405   40.   Inoue, F. *et al.* A systematic comparison reveals substantial differences in chromosomal

406        versus episomal encoding of enhancer activity. *Genome Res.* **27**, 38–52 (2017).

407

408

409 **Supplementary Methods**

410 *Cell culture*

411 We cultured K562 cells (ATCC) in IMDM (Gibco #12440) supplemented with 10% fetal bovine

412 serum (FBS) and 1% pen/strep and maintained in a humidified chamber at 37°C with 5% $CO_2$.

413 We cultured HepG2 cells (ATCC) in EMEM (ATCC #30-2003) supplemented with 10% FBS

414 and 1% pen/strep, maintained in a humidified chamber at 37°C with 5% $CO_2$.

415

416 *Generating an ORI-STARR-seq input plasmid library*

417 We sonicated human male genomic DNA (Promega #G1471) using a Covaris S220 sonicator

418 (duty factor – 5%; cycle per burst – 200; 40 sec) and ran it on a 0.8% agarose gel to size-select

419 500 bp fragments. After gel purification using a MinElute Gel Extraction kit (Qiagen), we end-

420 repaired, ligated custom adaptors, and PCR-amplified DNA fragments using Q5 Hot Start High-

421 Fidelity DNA polymerase (NEB) (98°C for 30 sec; 10 cycles of 98°C for 10 sec, 65°C for 30 sec,

422 and 72°C for 30 sec; 72°C for 2 min) to add homology arms for Gibson assembly cloning.

423 We used AgeI-HF (NEB) and SalI-HF (NEB) to linearize the hSTARR-seq_ORI plasmid (gift

424 from Alexander Stark; Addgene plasmid #99296) and cloned the PCR products into the vector

425 using Gibson Assembly Master Mix (NEB); we set up 60 replicate reactions to maintain

426 complexity. We purified the assembly reactions using SPRI beads (Beckman Coulter), dialyzed

427 them using Slide-A-Lyzer MINI dialysis devices (ThermoScientific), and concentrated them

428 using an Amicon Ultra-0.5 device (Amicon). We transformed the reaction into MegaX

429 DH10BTM T1 electrocompetent cells (Thermo Fisher Scientific) (with 25 replicate

430 transformations to maintain complexity) and let them grow in 12.5L LB-Amp medium until they

431 reached an optical density of ~1.0. We extracted the plasmids using a Plasmid Gigaprep Kit

20

432     (Qiagen) and dialyzed the plasmid prep using Slide-A-Lyzer MINI dialysis devices before

433     electroporation.

434

435     *Electroporation-mediated transfection of ORI-STARR-seq input plasmid library into K562 and*

436     *HepG2 cell lines*

437     We electroporated the ORI-STARR-seq library using an AgilePulse Max (Harvard Apparatus)

438     and generated two biological replicate for each cell line. For K562 cells, we electroporated 5.6

439     mg of input plasmid library into 700 million cells per biological replicate by delivering three 500

440     V pulses (1 ms duration with a 20 ms interval). For HepG2 cells, we electroporated 8 mg of input

441     plasmid library into one billion cells in one replicate, and 5.6 mg into 700 million cells in another

442     replicate by delivering three 300 V pulses (5 ms duration with a 20 ms interval).

443

444     *Generation of an Illumina sequencing library*

445     *Output RNA library*: We harvested cells 24 hr after electroporation, and extracted total RNA

446     using an RNeasy Maxi kit (Qiagen). We further isolated polyA-plus mRNA using Dynabeads®

447     Oligo (dT) kit (ThermoFisher Scientific), treated it with TURBO DNase (Invitrogen), and

448     purified the reaction using an RNeasy MinElute Kit (Qiagen). We synthesized cDNA using

449     SuperScript III (ThermoFisher Scientific) with a custom primer that specifically recognizes

450     mRNAs that had been transcribed from the ORI-STARR-seq library. After reverse transcription,

451     we treated the reactions with a cocktail of RNase A and RNase T1 (ThermoFisher Scientific).

452     We split cDNA samples into 160 replicate sub-reactions, and PCR-amplified each sub-reaction

453     with a primer with a unique index (helping to identify PCR duplicates) using Q5 Hot Start High-

454     Fidelity DNA polymerase (NEB) with the following program: 98°C for 30 s; cycles of 98°C for

21

455  10 s, 65°C for 30 s, 72°C for 30 s (until they reached mid-log amplification phase; we cycled 18

456  cycles for K562 Rep.1; 16 cycles for K562 Rep. 2; 18 cycles for HepG2 Rep. 1; and 15 cycles

457  for HepG2 Rep2); 72°C for 2 min). After PCR, we re-combined all sub-reactions into one and

458  purified it with Agencourt Beads. We generated 100 bp paired-end reads for each biological

459  replicate on an Illumina Hiseq4000 at the University of Chicago Genome Facility.

460  *Input DNA library*: We PCR-amplified a total of 200 ng of input plasmid library (in 16 replicate

461  reactions) using Q5 Hot Start High-Fidelity DNA polymerase (NEB) with the following

462  program: 98°C for 30 s; 4 cycles of 98°C for 10 s, 65°C for 30 s, and 72°C for 20 s; 8 cycles of

463  98°C for 10 s and 72°C for 50 s; 72°C for 2 min). After PCR, we combined all products into one

464  and purified it with Agencourt Beads. We generated 100 bp paired-end reads on an Illumina

465  Hiseq4000 at the University of Chicago Genome Facility.

466

467  *Sequencing and preprocessing*

468  For each of 160 replicates, paired-end sequencing reads were aligned to the human reference

469  genome hg38 using BWA-mem (v0.7.17). Alignments were filtered against unmapped,

470  secondary alignments, mapping quality score less than 30, and PCR duplicates using SAMtools

471  (v1.5) and Picard (v2.9.0). All of replicates were pooled and sorted for downstream analysis.

472

473  *Negative binomial distribution*

474  A negative binomial distribution, which arises from Gamma-Poisson mixture, can be

475  parametrized for y>=0 as follows.

476

$$Pr(Y = y_i | \mu_i, \theta) = f_Y(y_i; \mu_i, \theta) = \frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)} \cdot \left( \frac{\theta}{\theta + \mu_i} \right)^{\theta} \cdot \left( \frac{\mu_i}{\theta + \mu_i} \right)^{y_i}$$

477

478    Rearranging gives:

479

$$f_Y(y_i; \mu_i, \theta) = \frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)} \cdot \left(\frac{1}{1 + \frac{\mu_i}{\theta}}\right)^{\theta} \cdot \left(\frac{\frac{\mu_i}{\theta}}{1 + \frac{\mu_i}{\theta}}\right)^{y_i}$$

$$f_Y(y_i; \theta, \mu_i) = \frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)} \cdot \left(\frac{\mu_i}{\theta}\right)^{y_i} \left(\frac{1}{1 + \frac{\mu_i}{\theta}}\right)^{\theta + y_i}$$

$$f_Y(y_i; \theta, \mu_i) = \frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)} \cdot \left(\frac{\mu_i}{\theta}\right)^{y_i} \left(\frac{\theta}{\theta + \mu_i}\right)^{\theta + y_i}$$

$$f_Y(y_i; \theta, \mu_i) = \frac{\Gamma(y_i + \theta)}{\Gamma(y_i + 1) \cdot \Gamma(\theta)} \cdot \frac{\mu_i{}^{y_i} \theta^{\theta}}{(\theta + \mu_i)^{\theta + y_i}}$$

480

481    *Alternative parametrization of negative binomial regression using a rate model*

482    Alternative parametrization allows STARR-seq data to be modelled as a rate model. In contrast

483    to using input coverage as one of the covariates, we can consider it as "exposure" to output

484    coverage. This "trick" allows us to directly model the basal transcription rate (the ratio of RNA

485    to DNA) as a rate response variable. We defined the transcription rate (RNA to DNA ratio) as a

486    new variable, $\pi_i$.

487

$$\frac{y_i}{t_i} = \pi_i$$

488

23

489     If we assume the majority of genomic bins will have the basal transcription rate, we can model

490     the transcription rate at each $i$-th bin following the traditional negative binomial (NB2)

491     distribution.

492

$$\pi_i \sim NB\left(\frac{\mu_i}{t_i}, \theta\right)$$

493

494     The expected basal transcription, $E(\pi_i)$, becomes the mean incidence rate of $y_i$ per unit of

495     exposure, $t_i$.

496

$$E\left(\frac{y_i}{t_i}\right) = \frac{\mu_i}{t_i}$$

497

498     By normalizing $\mu_i$ by $t_i$, we are modeling a rate instead of a discrete count using the negative

499     binomial distribution. The regression term for the expected transcription rate can be expressed in

500     terms of a linear combination of explanatory variables, $j$ covariates ($\vec{x}$).

501

$$\ln\frac{\mu_i}{t_i} = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_j x_{ij}$$

502

503     Rearranging in terms of the expected value of $y$, or $\mu$, gives

504

$$\ln\mu_i - \ln t_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_j x_{ij}$$

$$\ln\mu_i = \ln t_i + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_j x_{ij}$$

$$\mu_i = \exp\left(\ln t_i + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_j x_{ij}\right)$$

24

505

506    The natural log of $t_i$ on the RHS ensures $\mu_i$ is normalized in the model, acting as an offset

507    variable. In STARRPeaker software, we allow users to optionally choose this alternative rate

508    model (implemented as "mode 2") instead of the default covariate model described in the main

509    text.

510

511    *BasicSTARRseq*

512    We used BasicSTARRseq R package version 1.10.0 downloaded from Bioconductor

513    (https://bioconductor.org/packages/release/bioc/html/BasicSTARRseq.html). We used default

514    setting as described in the software manual (minQuantile = 0.9, peakWidth = 500, maxPval =

515    0.001, deduplicate = TRUE, model = 1) to call peaks.

516

517    *MACS2*

518    We used MACS2 version 2.1.1 [23] at the recommended default setting, except for allowing

519    duplicates in read (--keep-dup all), since our STARR-seq dataset was multiplexed. We called

520    peaks with an FDR cutoff of 0.01, as recommended by the author of the software.

521
522    **Supplementary Tables**

523    Table S1 contains significant peaks called by STARRPeaker.

524    Table S2 contains various statistics from comparing STARRPeaker peaks to peaks called by

525    BasicSTARRseq and MACS2.

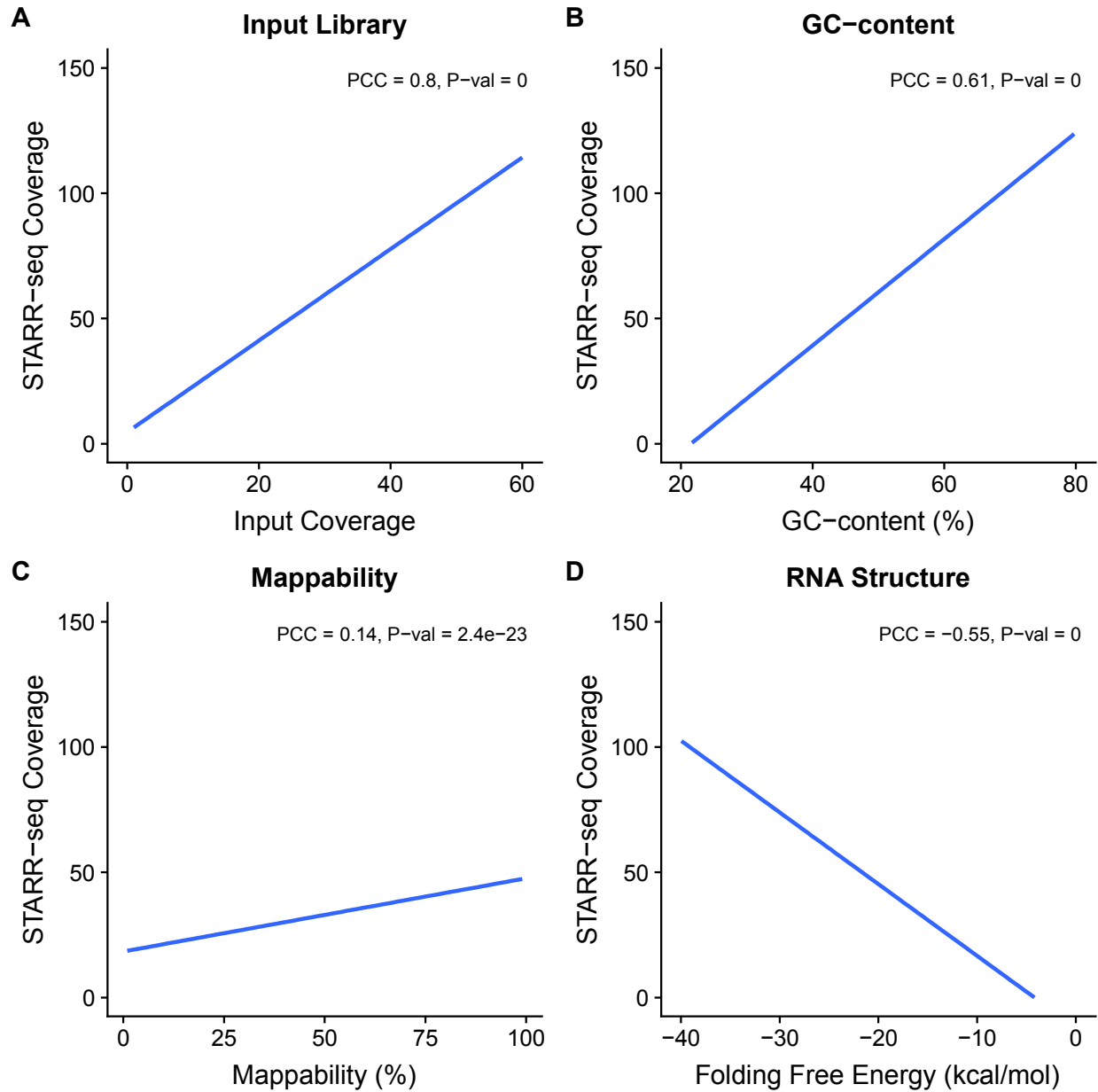526    Table S3 contains list of data sources and accession number used for the analysis.

527

528  **Figures**



STARR−Peaker (Coverage using fragment center)
Others (Coverage using read star t)

529

530  **Figure 1** *Comparison of STARR-seq coverage calculated using fragment center to*
531  *using read start position. (A)-(D) shows examples drawn from K562 STARR-seq data.*
532  *Triangle indicates the summit of coverage. Read depth was normalized, since 2 paired*
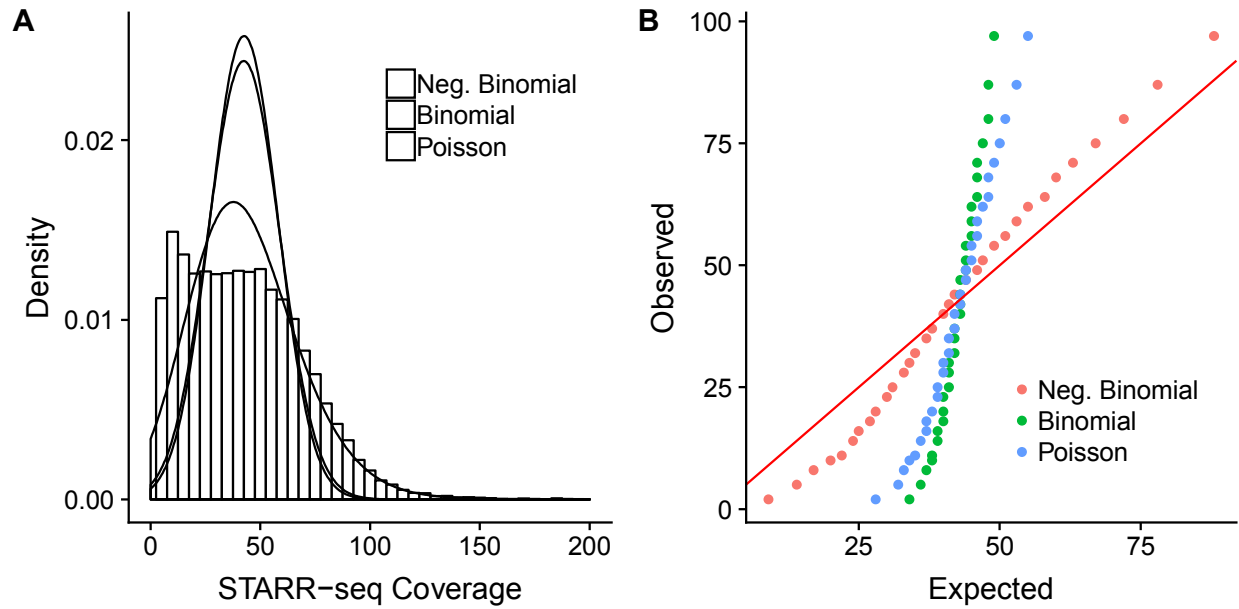533  *reads correspond to 1 fragment.*

534

**Figure 2** *Confounding factors in the STARR-seq assay. STARR-seq output and input coverages are significantly correlated with (A) input coverage (B) GC-content (C) mappability, and (D) RNA structure folding. PCC: Pearson Correlation Coefficient. Plots were from a sampling of 5,000 random genomic bins.*

541

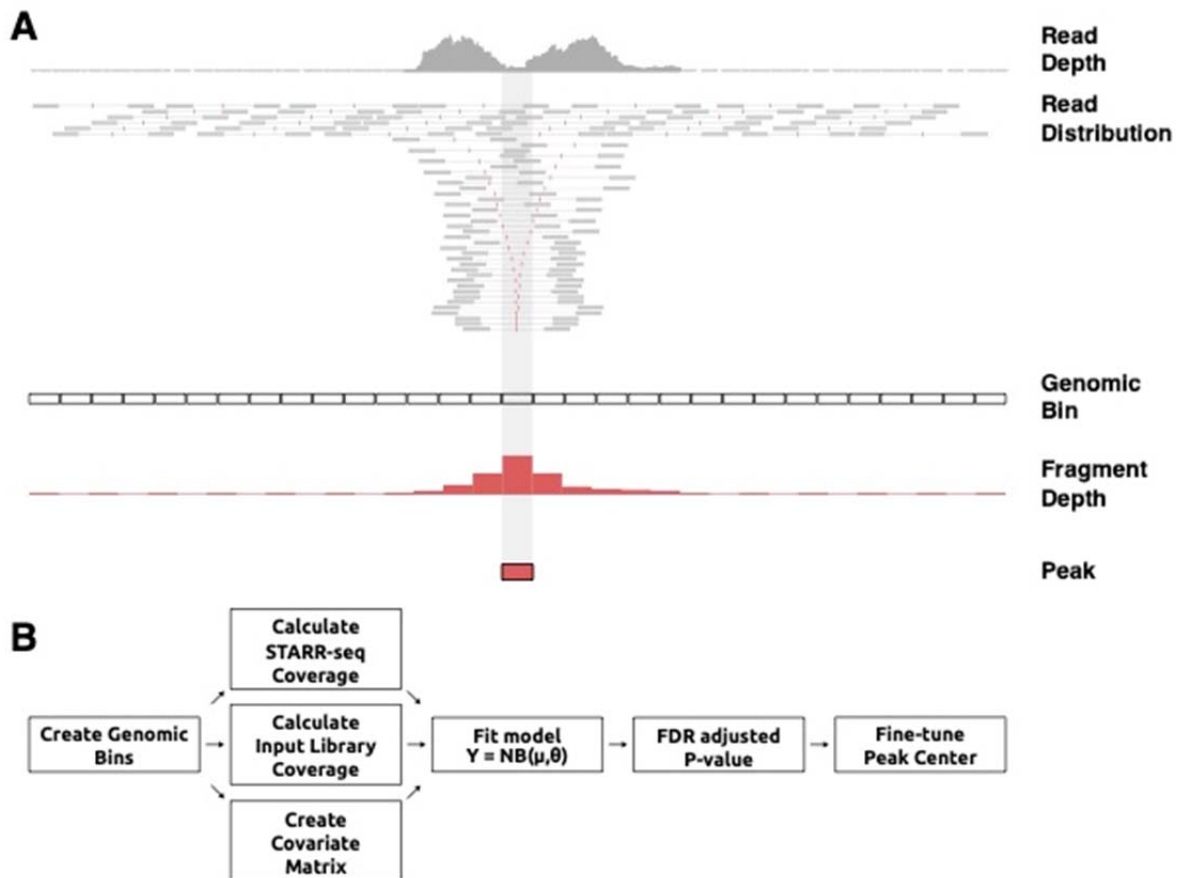

**Figure 3** *STARR-seq coverage is fitted against simulated coverage using three distribution models; negative binomial, binomial, and Poisson. (A) Density histogram of simulated distribution against STARR-seq coverage. (B) Q-Q plot of simulated distribution against STARR-seq coverage. The red solid line represents where the observed count equals the expected count.*
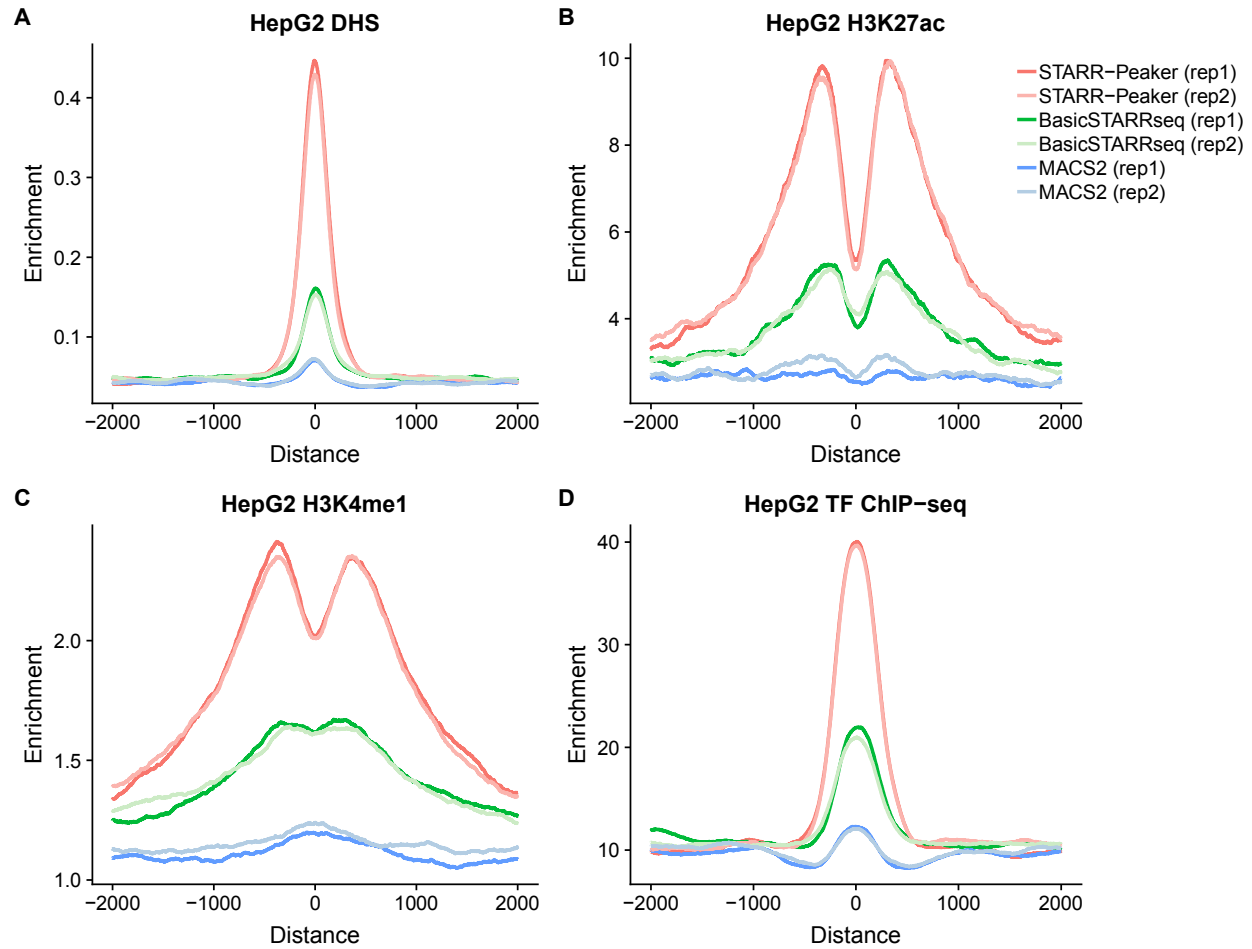
547

548

*Figure 4* *Overview of STARRPeaker peak-calling scheme. (A) In contrast to using read depth (grey), fragment depth (red) offers more precise and sharper STARR-seq coverage. Fragment inserts are directly inferred from properly paired-reads. (B) Workflow of STARRPeaker describing how coverage is calculated for each genomic bin and modelled using negative binomial regression model. The analysis pipeline can largely be divided into four steps: (1) Binning the genome (2) Calculating coverage and computing covariate matrix (3) Fitting the STARR-seq data to the NB regression model (4) Peak calling, multiple hypothesis testing correction, and adjustment of the center of peaks*
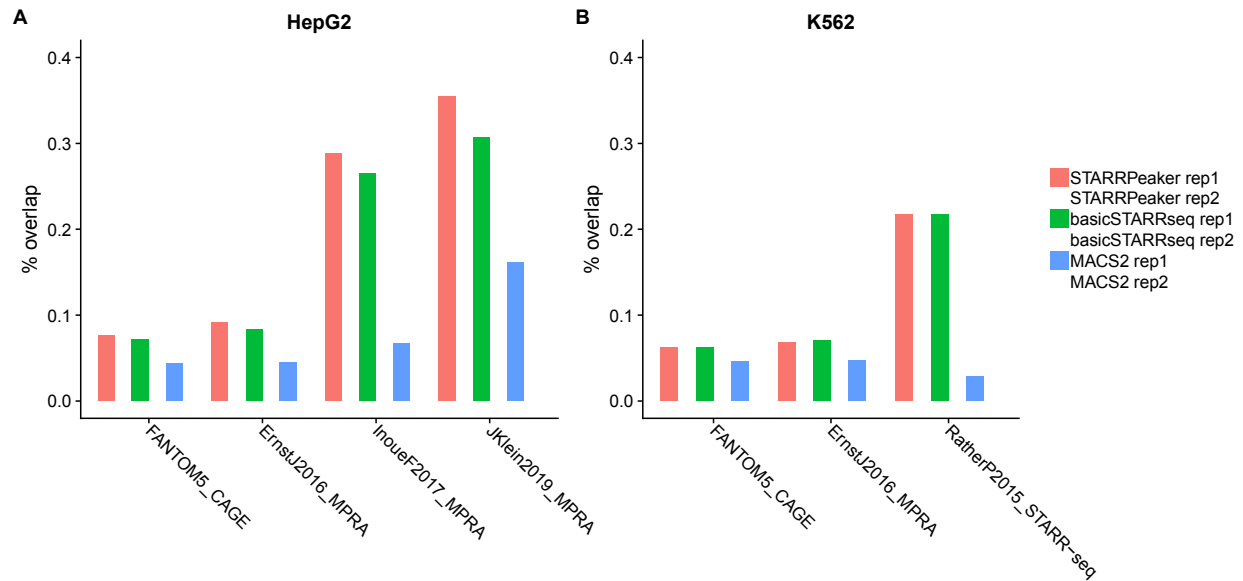
557

29

558

*Figure 5 Enrichment of epigenetic signals around peaks. All peaks were centered at the summit, uniformly thresholded using P-value < 0.001, and 10,000 peaks were randomly selected. Aggregated read depth at 2000 bp upstream and downstream were plotted for (A) DNase-seq (B) H3K27ac (C) H3K4me1 (D) Aggregated TF ChIP-seq profile. For TF ChIP-seq, high enrichment indicates TF binding hotspots*

564

565

**Figure 6** *Comparison of peaks using external dataset. Peaks identified from STARRPeaker as well as BasicSTARRseq and MACS2 were compared against published dataset. For a fair comparison, 100,000 peaks were randomly drawn from peaks identified by each peak caller using the recommended settings, and the fraction of overlap was computed for each replicate. We considered it as an overlap when at least 50% of peaks intersected each other.*