

On accounting for sequence-specific bias in genome-wide chromatin accessibility experiments: recent advances and contradictions

Pedro Madrigal^{1, 2*}

¹Wellcome Trust Sanger Institute, United Kingdom, ²Department of Surgery, University of Cambridge, United Kingdom

Submitted to Journal:
Frontiers in Bioengineering and Biotechnology

Specialty Section:
Bioinformatics and Computational Biology

ISSN:
2296-4185

Article type:
Opinion Article

Received on:
14 Jun 2015

Accepted on:
07 Sep 2015

Frontiers website link:
www.frontiersin.org

Citation:
Madrigal P(2015) On accounting for sequence-specific bias in genome-wide chromatin accessibility experiments: recent advances and contradictions. *Front. Bioeng. Biotechnol.* 3:144. doi:10.3389/fbioe.2015.00144

Copyright statement:
© 2015 Madrigal. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.

1
2
3
4 On accounting for sequence-specific bias in genome-wide chromatin
5 accessibility experiments: recent advances and contradictions
6
7

8 Pedro Madrigal^{1,2}
9

10 ¹ Wellcome Trust Sanger Institute, Hinxton, Cambridge, UK

11 ² Department of Surgery, University of Cambridge, UK
12
13
14
15

16 **Correspondence:**

17
18 Wellcome Trust Sanger Institute
19 Wellcome Trust genome Campus, Hinxton
20 Cambridge, CB10 1SA, UK
21 pm12@sanger.ac.uk
22
23

24 **Keywords: (1-8)**

25 Next generation sequencing, DNase-seq, ATAC-seq, chromatin accessibility, footprinting, sequence bias, ChIP-
26 exo
27
28
29
30
31
32
33
34
35
36
37
38
39

40 Number of words: 1459

41 Number of figures: 1
42
43
44
45
46

Next-generation sequencing for chromatin biology

Uncovering the protein-DNA interactions involved in cell fate, development, and disease in a time- and cell-specific manner is a fundamental goal of molecular biology. The advent of the sequencing technologies has opened a new genomic era, uncovering the information encoded in genomes, epigenomes, and transcriptomes (McPherson, 2014). For example, the popular ChIP-based techniques ChIP-seq (Johnson et al., 2007; Robertson et al., 2007) and ChIP-exo (Rhee and Pugh, 2011) are widely used to detect transcription factor (TF) binding sites using an antibody against a single protein of interest (Mahony and Pugh, 2015). Alternative protocols assaying the chromatin landscape, such as those based on digestion by DNase I enzyme (DNase-seq), micrococcal nuclease (MNase-seq), and Tn5 transposase attack (ATAC-seq), enable the identification of DNA-binding protein footprints of many TFs in a single experiment (Tsompana and Buck, 2014). Time series experiments might be required for the identification of those TFs catalogued as pioneer factors, allowing their effects on chromatin to be investigated (Zaret and Carroll, 2011; Sherwood et al., 2014; Pajoro et al., 2014).

Despite the initial promise of detecting the majority of TFs in one assay, DNA sequence-specific biases, together with TF-dependent binding kinetics, have been recently pinpointed as major confounding factors in DNase-seq experiments (He et al., 2014; Koohy et al., 2013; Raj and McVicker, 2014; Rusk, 2014; Sung et al., 2014). These influencing factors were not considered by any of the previous computational approaches for the analysis of next-generation sequencing chromatin accessibility data (Madrigal and Krajewski, 2012); neither those strategies based on TF-generic DNase signature nor those based on TF-specific DNase signature (Luo and Hartemink, 2013).

Alleviating sequence-specific biases in DNase-seq

To partly address these challenges, four recent approaches have been published that model, predict, or explain DNase I sequence specificity in order to improve the detection of TF occupancy events at high resolution (digital genomic footprinting). The first method, FootprintMixture, uses a multinomial mixture model in which one mixture models the footprint component, and the other the background component taking into account the sequence bias (Yardimci et al., 2014). The background can be either uniform or derived from naked DNA measurements - this is the main difference with respect to the footprint component in CENTIPEDE (Pique-Regi et al., 2011), which assumes a uniform background. Alternatively, more than 2 components may be set to detect variability in the footprint model. Thus, the cleavage signature (number of DNase I cuts that map to each nucleotide) is used in a multinomial mixture model to classify candidate sites as either 'bound' or 'unbound' aided by 6-mer DNase sequence bias cleavage frequencies (Yardimci et al., 2014). Remarkably, the authors found that sequence bias is DNase-seq protocol-specific. They also found that the signature of a footprint could be formed by a mixture of DNase digestion profiles identified by unsupervised *k*-means clustering, in agreement with the observations found in an earlier study (Tewari et al., 2012). For transcription factors CTCF and ZNF143, variants of the consensus sequence motif associated to different footprint shapes were observed.

In the second, the DNase2TF algorithm is able to correct dinucleotide bias, detecting footprints with accuracy better or comparable to existing approaches (Sung et al., 2014). Furthermore, Sung et al. (2014) were able to predict DNase signatures using solely tetranucleotide frequency information. Although this 4-nucleotide region has the highest information content, Koohy et al. (2013) and Lazarovici et al. (2013) demonstrated information beyond a context longer than 4 nucleotides. Consequently, using naked (deproteinized) DNA control datasets specific to a protocol and an enzyme, as well as high sequencing depth (Hesselberth et al., 2009), are now suggested recommendations for DNase-seq experiments aiming to detect footprints (Meyer and Liu, 2014).

A third approach, an improved version of HINT (HMM-based identification of transcription factor footprints (Gusmao et al., 2014)), named as HINT-BC/ HINT-BCN (Bias Correction based on hypersensitivity sites/ Bias Correction based on Naked DNase-seq) includes *k*-mer based bias correction in DNase-seq data as in He et al. (2014), leading to substantial changes in the average DNase I cleavage patterns surrounding the TFs. These changes result beneficial to footprinting method accuracy (personal communication with the author). Contradictorily, a fourth study using DNase-seq has shown that bias correction does not significantly improve the accuracy of TF binding identification (Kähärä and Lähdesmäki, 2015). In addition, this study poses a second counterintuitive idea in the field: accuracy saturates at a modest sequencing depth (30-60 million reads), and only a few TFs present improvement at deeper sequencing.

ATAC-seq shows sequence cleavage bias

It is unknown if ATAC-seq derived footprints are factor-dependent or affected by Tn5 cleavage preferences (Tsompana and Buck, 2014). As expected, bioinformatic analysis of chromosome 22 in the published human datasets for 50,000 cells reveals sequence biases in ATAC-seq experiments (Buenrostro et al., 2013) (Figure 1), similar to those found by Koohy et al. (2013) in DNase-seq. As ATAC-seq might replace DNase-seq in the foreseeable future due to its cost and time efficiencies, and because it simultaneously allows the identification of nucleosome positions (Buenrostro et al., 2013), new computational models are necessary to evaluate intrinsic confounding factors in ATAC-seq.

A novel approach, msCentipede (Raj et al., 2014), has extended CENTIPEDE (Pique-Regi et al., 2011) from a multinomial model to a hierarchical multiscale model. It has been evaluated on 'single-hit' UW DNase-seq (Hesselberth et al., 2009) and on paired-end (PE) ATAC-seq data. Surprisingly, the 'flexible model' for background DNase I cleavage rate (msCentipede-flexbg) shows very little improvement for a broad range of factors when taking into account naked DNA information from Lazarovici et al. datasets (2013). This finding clearly contradicts those of He et al. (2014) and Sung et al. (2014). In msCentipede, the footprint signature (or cleavage profile) pattern within a factor-bound motif instance was therefore found to be informative when increasing the sensitivity and specificity of the TF binding site prediction. Raj et al. (2014) suggest that this might be explained by the different range of read count data between the matched consensus sequence of the candidate site/motif (10-30 bp) and the data matrix used typically by the software packages (larger sequence window, around 100-150 bp extension at each flank of the motif), which can mask the effects produced by not accounting for sequence biases within the core motif.

Are current benchmarks adequate to evaluate bias-corrected DNase-seq data?

So far, a footprint of a TF therefore might be either detectable (and better detectable when accounting, or not, for influencing factors), or undetectable. In many studies, both problems are convoluted and addressed using the same "gold standard" datasets, such as ChIP-seq, which do not have nucleotide-level resolution. Hence, on these methods and gold standards no reproducible improvements can be seen. This was already noted in Cuellar-Partida et al. (2012), when it was showed that simply scanning for position weight matrices in DNase I hypersensitive sites (DHSs) had the same power as CENTIPEDE. These issues also complicate data integration with TF ChIP-seq, as peaks without a footprint in DNase-seq/ATAC-seq, considered weak/indirect binding or false positives (ChIP artifacts), might instead be explained by a class of TFs with rapid kinetics. And vice versa, DNase I cleavage patterns located within 'ChIP-seq unbound' sites - noted previously, e.g., in the MILLIPEDE framework, especially in yeast (Luo and Hartemink, 2013)- could support the hypothesis of footprint shape dominated by DNA sequence specificities.

135 **Future directions**

136 There is room for improvement in current methodologies by making use of the sequence-specificity of each
137 enzyme/assay, including ATAC-seq, but there is no clear consensus in its importance for digital genomic
138 footprinting. This situation is not exclusive for genome-wide chromatin accessibility experiments: modelling the
139 sequence-specific lambda exonuclease bias in ChIP-exo did not significantly increase the identification of TF
140 binding sites (Wang et al., 2014). Similarly, there is no clear consensus if footprint signatures at the core motif,
141 whether they are unique or not for an individual factor, are really important for footprint identification.
142 Establishing better benchmarks to compare performance of the algorithms across different protocols is a
143 fundamental task. These benchmarks could be based on “differential footprints” (sites within DHSs that are
144 bound by a factor in one condition but not the other) as a more appropriate metric to evaluate footprint
145 identification performance instead of using ChIP-seq data (Yardimci et al., 2014). In addition, are DNase-seq
146 software tools equally applicable to ATAC-seq without modification? If enzyme-specific biases are taken into
147 account in a comparable experimental set-up, will DNase-seq and ATAC-seq report the same footprints for an
148 identical sample using same algorithm parameters? This is unlikely, based on a previous comparison between
149 open chromatin DNase I hypersensitive sites and FAIRE sites, which revealed unique regions produced in each
150 assay (Song et al., 2011). It has been also proposed that performing, and combining, experiments with different
151 nucleases can be an alternative to mitigate biases (He et al., 2014; Mahony and Pugh, 2015).
152 A greater challenge is dealing with proteins with very short residency time in the DNA, as they produce mostly
153 negligible footprints (Rusk, 2014; Sung et al., 2014). Optimizing and implementing new methods is necessary in
154 order to enable biological insights that current methods cannot reveal.

155

156

157

158

159 **DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT**

160 The author declares that the research was conducted in the absence of any commercial or financial relationships
161 that could be construed as a potential conflict of interest.

162

163

164 **FIGURES**

165 **Figure 1.** Tn5 transposase shows sequence cleavage bias. Data represented correspond to read-start sites in
166 reads aligned to forward and reverse strands in chromosome 22 in four ATAC-seq replicates (50k cells per
167 replicate) reported in Buenrostro et al. (2013). 50 bp PE reads were pre-processed with Trimmomatic v0.32
168 under default parameters, and then aligned to hg19 using BWA v0.7.4-r385 (Bolger et al., 2014; Li and Durbin,
169 2010). Sequence logos were generated using WebLogo (Crooks et al., 2004). Y-axis: 0.0 - 0.4 bits.

REFERENCES

- Bolger, A. M., Lohse, M., and Usadel, B. (2014), Trimmomatic: a flexible trimmer for Illumina sequence data, *Bioinformatics*, 30, 15, 2114–2120
- Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y., and Greenleaf, W. J. (2013), Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position, *Nat. Methods*, 10, 12, 1213–1218
- Crooks, G. E., Hon, G., Chandonia, J. M., and Brenner, S. E. (2004), WebLogo: a sequence logo generator, *Genome Res.*, 14, 6, 1188–1190
- Cuellar-Partida G, Buske FA, McLeay RC, Whittington T, Noble WS, Bailey TL. (2012) Epigenetic priors for identifying active transcription factor binding sites. *Bioinformatics*, 28, 56-62
- Gusmao, E. G., Dieterich, C., Zenke, M., Costa, I. G. (2014) Detection of active transcription factor binding sites with the combination of DNase hypersensitivity and histone modifications, *Bioinformatics*, 30, 3143-3151
- He, H. H., Meyer, C. A., Hu, S. S., Chen, M. W., Zang, C., Liu, Y., et al. (2014), Refined DNase-seq protocol and data analysis reveals intrinsic bias in transcription factor footprint identification, *Nat. Methods*, 11, 1, 73–78
- Hesselberth, J. R., Chen, X., Zhang, Z., Sabo, P. J., Sandstrom, R., Reynolds, A. P., et al. (2009), Global mapping of protein-DNA interactions in vivo by digital genomic footprinting, *Nat. Methods*, 6, 4, 283–289
- Johnson, D. S., Mortazavi, A., Myers, R. M., and Wold, B. (2007), Genome-wide mapping of in vivo protein-DNA interactions, *Science*, 316, 5830, 1497–1502
- Kähärä J, Lähdesmäki H (2015) BinDNase: a discriminatory approach for transcription factor binding prediction using DNase I hypersensitivity data. *Bioinformatics*. 2015 May 7. pii: btv294.
- Koohy, H., Down, T. A., and Hubbard, T. J. (2013), Chromatin accessibility data sets show bias due to sequence specificity of the DNase I enzyme, *PLoS ONE*, 8, 7, e69853
- Lazarovici, A., Zhou, T., Shafer, A., Dantas Machado, A. C., Riley, T. R., Sandstrom, R., et al. (2013), Probing DNA shape and methylation state on a genomic scale with DNase I, *Proc. Natl. Acad. Sci. U.S.A.*, 110, 16, 6376–6381
- Li, H. and Durbin, R. (2010), Fast and accurate long-read alignment with Burrows-Wheeler transform, *Bioinformatics*, 26, 5, 589–595
- Luo, K. and Hartemink, A. J. (2013), Using DNase digestion data to accurately identify transcription factor binding sites, *Pac Symp Biocomput*, 80–91
- Madrigal, P. and Krajewski, P. (2012), Current bioinformatic approaches to identify DNase I hypersensitive sites and genomic footprints from DNase-seq data, *Front Genet*, 3, 230
- Mahony, S. and Pugh, B. F. (2015) Protein-DNA binding in high-resolution, *Crit Rev Biochem Mol Biol*, Jun 3:1-15.
- McPherson, J.D. (2014) A defining decade in DNA sequencing. *Nat Methods*, 11,1003-1005.
- Meyer, C. A. and Liu, X. S. (2014), Identifying and mitigating bias in next-generation sequencing methods for chromatin biology, *Nat. Rev. Genet.*, 15, 11, 709–721
- Pajoro, A., Madrigal, P., Muino, J. M., Matus, J. T., Jin, J., Mecchia, M. A., et al. (2014), Dynamics of chromatin accessibility and gene regulation by MADS-domain transcription factors in flower development, *Genome Biol.*, 15, 3, R41
- Pique-Regi, R., Degner, J. F., Pai, A. A., Gaffney, D. J., Gilad, Y., and Pritchard, J. K. (2011), Accurate inference of transcription factor binding from DNA sequence and chromatin accessibility data, *Genome Res.*, 21, 3, 447–455
- Raj, A. and McVicker, G. (2014), The genome shows its sensitive side, *Nat. Methods*, 11, 1, 39–40
- Raj, A., Shim, H., Gilad, Y., Pritchard, J. K., and Stephens, M. (2014), mscentipede: Modeling heterogeneity across genomic sites improves accuracy in the inference of transcription factor binding, *bioRxiv*,

- 217 Rhee, H. S. and Pugh, B. F. (2011), Comprehensive genome-wide protein-DNA interactions detected at single-
218 nucleotide resolution, *Cell*, 147, 6, 1408–1419
- 219 Robertson, G., Hirst, M., Bainbridge, M., Bilenky, M., Zhao, Y., Zeng, T., et al. (2007), Genome-wide profiles
220 of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing, *Nat.*
221 *Methods*, 4, 8, 651–657
- 222 Rusk, N. (2014), Transcription factors without footprints, *Nat. Methods*, 11, 10, 988–989
- 223 Sherwood, R. I., Hashimoto, T., O'Donnell, C. W., Lewis, S., Barkal, A. A., van Hoff, J. P., et al. (2014),
224 Discovery of directional and nondirectional pioneer transcription factors by modeling DNase profile magnitude
225 and shape, *Nat. Biotechnol.*, 32, 2, 171–178
- 226 Song, L., Zhang, Z., Grasfeder, L. L., Boyle, A. P., Giresi, P. G., Lee, B. K., et al. (2011), Open chromatin
227 defined by DNaseI and FAIRE identifies regulatory elements that shape cell-type identity, *Genome Res.*, 21, 10,
228 1757–1767
- 229 Sung, M. H., Guertin, M. J., Baek, S., and Hager, G. L. (2014), DNase footprint signatures are dictated by factor
230 dynamics and DNA sequence, *Mol. Cell*, 56, 2, 275–285
- 231 Tewari, A. K., Yardimci, G. G., Shibata, Y., Sheffield, N. C., Song, L., Taylor, B. S., et al. (2012), Chromatin
232 accessibility reveals insights into androgen receptor activation and transcriptional specificity, *Genome Biol.*, 13,
233 10, R88
- 234 Tsompana, M. and Buck, M. J. (2014), Chromatin accessibility: a window into the genome, *Epigenetics*
235 *Chromatin*, 7, 1, 33
- 236 Wang, L., Chen, J., Wang, C., Uusküla-Reimand, L., Chen, K., Medina-Rivera, A. et al. (2014), MACE: model
237 based analysis of ChIP-exo. *Nucleic Acids Res.* 42(20), e156.
- 238 Yardimci, G. G., Frank, C. L., Crawford, G. E., and Ohler, U. (2014), Explicit DNase sequence bias modeling
239 enables high-resolution transcription factor footprint detection, *Nucleic Acids Res.*, 42, 19, 11865–11878
- 240 Zaret, K. S. and Carroll, J. S. (2011), Pioneer transcription factors: establishing competence for gene expression,
241 *Genes Dev.*, 25, 21, 2227–2241

Figure 1.JPEG

