

Analysis of computational footprinting methods for DNase sequencing experiments

Eduardo G. Gusmao^{1,2}, Manuel Allhoff^{1,3}, Martin Zenke^{1,2}, Ivan G. Costa^{1,2,3,*}.

¹ IZKF Computational Biology Research Group, RWTH Aachen University Medical School, Aachen, Germany.

² Department of Cell Biology, Institute of Biomedical Engineering, RWTH Aachen University Medical School, Aachen, Germany.

³ Aachen Institute for Advanced Study in Computational Engineering Science (AICES), RWTH Aachen University, Aachen, Germany.

* e-mail: ivan.costa@rwth-aachen.de

Final Version: Gusmao et al, Nature Methods, 13, 303–309, doi:10.1038/nmeth.3772
<http://www.nature.com/nmeth/journal/v13/n4/pdf/nmeth.3772.pdf>

Abstract: DNase-seq allows a nucleotide-level identification of transcription factor binding sites based on the computational search of footprint-like DNase I cleavage patterns on the DNA. Frequently, in high-throughput methods, experimental artifacts like DNase I cleavage bias impact the computational analysis of DNase-seq experiments. Here we performed a comprehensive and systematic study on the performance of computational footprinting methods. We evaluated 10 footprinting methods on a panel of DNase-seq experiments for their ability to recover cell-specific transcription factor binding sites. We show that three methods: HINT, DNase2TF and PIQ consistently outperform other evaluated methods. We demonstrate that correcting the DNase-seq signal for experimental artifacts significantly improves accuracy of computational footprints. We also propose a score to detect footprints arising from transcription factors with potentially short residence time.

Next-generation sequencing (NGS) combined with genome-wide mapping techniques, such as DNase-seq, contributed greatly to our understanding of gene regulation and chromatin dynamics^{1,2,3}. DNase-seq allows a nucleotide-level identification of transcription factor binding sites (TFBSs). This can be performed by the computational search of footprint-like regions with low number of DNase I cuts surrounded by regions with high number of cuts^{2,3}. A number of computational footprinting methods have been proposed in the past years⁴⁻¹³. Among other applications, these methods allow the delineation of the human regulatory lexicon with millions of TFBSs over distinct cell types⁴, the detection of uncharacterized transcription factor (TF) motifs indicating putative regulatory elements⁴ and the study of conservation of regulatory regions across different species¹⁴.

NGS-based data are significantly affected by artifacts, which are inherent to the experimental protocols used^{15,16,17}. An example is the DNase I sequence cleavage bias, which is due to DNase I having different binding affinities towards specific DNA sequences. He et al.¹⁵ showed that sequence cleavage bias around TFBSs strongly affects the performance of a computational footprinting method^{4,15} (footprint score; FS) in a TF-specific manner. They also indicated several TFs, such as nuclear receptors and *de novo* motifs found via computational footprinting⁴, where the DNase-seq profile resembles their sequence cleavage bias estimate. Furthermore, they indicated that ranking putative TFBS by the number of DNase-seq reads around putative TFBSs^{10,15} (tag count; TC) outperforms the ranking by FS. Another experimental aspect affecting the computational analysis of DNase-seq is the residence time of TF binding. Sung et al.⁷ showed that short-lived TFs display a lower DNase I cleavage protection pattern, i.e. low number of DNase-seq reads surrounding the footprint. Moreover, they also noticed that nuclear receptors have DNase-seq profiles resembling their DNase I sequence cleavage bias estimates. While both studies^{7,15} show the challenges imposed by cleavage bias and residence time, there have been a few attempts^{7,12,15} to address these computationally.

There is no well-defined gold standard for the evaluation of footprinting methods. All work so far has used ChIP-seq of TFs in conjunction with motif-based predictions as ground truth. In short, motif-predicted binding

48 sites (MPBSs) supported by ChIP-seq peaks are positive examples (true TFBSs), while MPBSs without
49 ChIP-seq support are negative examples (false TFBSs)¹⁰. This evaluation requires TF ChIP-seq experiments
50 to be carried out on the very same cells as the DNase-seq experiment and has a few caveats. First, TF
51 ChIP-seq peaks are also observed in indirect binding events^{4,7,12,18}. Second, they have a lower spatial
52 resolution than DNase-seq. Therefore, false TFBSs might be regarded as true TFBSs by proximity to a real
53 TFBS of a distinct TF^{15,17}. Recently, Yardımcı et al.¹² indicated that footprint quality scores, as measured by
54 the footprint likelihood ratio (FLR), were significantly higher in cells where the TF was expressed. This
55 observation indicates that comparing changes in expression and quality of footprints in a pairs of cells could
56 provide an alternative footprint evaluation measure. Finally, with the exception of a few studies^{8,11,12,13},
57 comparative analyses evaluating footprinting methods were based on ChIP-seq of few (<12) TFs and with
58 the exception of Gusmao et al.⁸, a maximum of four competing methods were evaluated. Despite the
59 importance of method evaluation¹⁹, there is a clear lack of benchmark data, evaluation standards and studies
60 performing a comprehensive analysis of computational footprinting methods.

61 We evaluated 10 computational footprinting methods: Neph⁴, Boyle⁵, Wellington⁶, DNase2TF⁷, HINT⁸,
62 Centipede⁹, Cuellar¹⁰, PIQ¹¹, FLR¹² and BinDNase¹³. In a “ChIP-seq based approach” they are evaluated in
63 their accuracy to recover TFBSs supported by 88 ChIP-seq TF experiments of two cell types (H1-hESC and
64 K562) with the area under the receiver operating characteristic curves (AUC) and precision-recall curves
65 (AUPR). We also propose the “FLR-Exp” methodology, which associates the FLR¹² scores for footprints in
66 cell type pairs with the fold change expression of the TFs associated to the footprints. This analysis is based
67 on the comparison of footprints and expression of 143 TFs in H1-hESC, K562 and GM12878 cells. We also
68 evaluate approaches for ranking footprints, strategies for dealing with DNase-seq experimental artifacts and
69 the effect of TF residence time on footprint predictions.

70 RESULTS

71 Computational genomic footprinting methods

72 Computational footprinting methods can be broadly categorized in segmentation (SEG)⁴⁻⁸ and site-centric
73 (SC) methods⁹⁻¹³. Several segmentation methods use window search to scan DNase-seq genomic profiles
74 with a footprint-like shape – short regions with low DNase-seq digestion between short regions with high
75 DNase-seq digestion (Neph⁴, Wellington⁶ and DNase2TF⁷). Another family of segmentation methods are
76 based on hidden Markov models (HMMs), in which the hidden states model distinct levels of DNase-seq
77 cleavage activity around footprints (Boyle⁵ and HINT⁸). Site-centric methods analyze DNase-seq profiles
78 around MPBSs and classify these sites as being either bound or unbound. Most site-centric methods are
79 based on unsupervised statistical methods like mixture models (FLR¹²), Bayesian mixture models
80 (Centipede⁹) and combination of Gaussian process (GP) and expectation propagation (PIQ¹¹). An alternative
81 site-centric approach is proposed by Cuellar¹⁰, which uses DNase-seq profiles as prior distribution for the
82 detection of MPBSs. BinDNase is a supervised site-centric method based on logistic regression¹³. We also
83 evaluate simple statistics as baseline methods: ranking MPBSs by position weight matrix (PWM-Rank) bit-
84 score¹⁰, by ratio of the number of DNase-seq reads inside and around a MPBS (FS-Rank)^{4,15} and by number
85 of DNase-seq reads around a MPBS (TC-Rank)^{10,15}.

86 There are several other relevant characteristics for computational footprinting methods. A few methods allow
87 the inclusion of additional genomic and/or experimental evidence like conservation scores⁹, distance to
88 transcription start sites⁹ and histone modifications⁸⁻¹⁰. Only PIQ¹¹ supports the analysis of several DNase-seq
89 data sets, i.e. experiments with replicates or time series. Another important feature is the correction of
90 DNase-seq experimental artifacts, which is only supported by DNase2TF⁷, HINT⁸ variants (HINT-BC and
91 HINT-BCN) and FLR⁹. While HINT-BC, HINT-BCN and DNase2TF use experimental bias statistics to pre-
92 process DNase-seq profiles; FLR builds a “cleavage bias” model within their mixture model in a TF-specific
93 manner. Most methods use base pair DNase-seq resolution as primary input^{4-9,11-13}. One exception is
94 Cuellar¹⁰, which is based on smoothed DNase-seq signals of windows with 150 bps. Smoothing of base pair
95 resolution profiles is performed by PIQ via the use of GP models¹¹. BinDNase uses a greedy backward
96 feature selection approach, which merges read counts of neighboring genomic positions¹³. Footprinting

97 methods also provide statistics to rank footprint predictions. Wellington⁶ and DNase2TF⁷ use read count
98 statistics to provide *p*-values for each footprint. Several site-centric approaches provide either probabilities
99 (BinDNase¹³, Centipede⁹ and PIQ¹¹) or log-odds scores (FLR¹²) of footprints. Other methods use statistics
100 such as FS (Neph⁴), PWM (Cuellar¹⁰) scores or TC (HINT⁸), to rank predicted footprints.

101 The availability, usability and scalability of software tools implementing the methods are also important
102 features. Neph⁴, HINT⁸, PIQ¹¹ and Wellington⁶ provide tutorials and software to run experiments with few
103 command line calls. Of those, only HINT⁸, PIQ¹¹ and Wellington⁶ natively support standard genomic formats
104 as input. Site centric methods Cuellar¹⁰, BinDNase¹³, Centipede⁹ and FLR¹² require a single execution and
105 input data per TF and cell, while segmentation methods require an execution per cell only. These site centric
106 methods have computational demands 5 times (FLR and Cuellar) to 50 times (BinDNase and Centipede)
107 higher than the slowest segmentation method (Wellington) on our analysis (**Supplementary Table 1**). The
108 main method features are summarized in **Table 1** and described in the **Online Methods**.

109 **Association of TF expression with footprint quality**

110 Yardımcı et al. indicated that the FLR of candidate footprints are significantly higher in cells where the TF is
111 being expressed¹². We expand this idea by evaluating if differences in FLR score distribution of footprints
112 overlapping with MPBSs on a pair of cell types are proportional to differences in the expression of the
113 respective TFs (**Fig. 1a**). We observed high average correlation values for the majority of evaluated methods
114 ($r = 0.79$) and extremely high correlation values ($r > 0.9$) for top performing methods on comparisons
115 between pairs of cell types H1-hESC, K562 and GM12878 (**Fig. 1b**; **Supplementary Fig. 1**). We also
116 evaluated the use of the TC and FS metrics as quality scores instead of FLR. They had lower average
117 correlation values (TC $r = 0.35$ and FS $r = 0.73$; **Supplementary Fig. 2**). We opt, therefore, to use the FLR
118 as quality measure for footprints for this evaluation procedure. The correlation between FLR score difference
119 and expression fold change, which we refer to as “FLR-Exp”, will be used to rank footprinting methods.
120 Highest values indicate best performance. The FLR-Exp evaluation methodology only requires expression
121 data and is therefore more generally applicable than TF ChIP-seq based evaluation. However, differently
122 from the TF ChIP-seq evaluation, the FLR-Exp approach cannot evaluate footprint predictions of individual
123 TFs.

124 **Impact of experimental artifacts**

125 To understand the nature of artifacts on DNase-seq experiments, we analyzed the sequence bias estimates
126 on all 61 ENCODE Tier 1 and 2 DNase-seq data sets (**Supplementary Table 2**). These experiments include
127 two main DNase-seq protocols, which differ on the number of DNase I digestion events necessary to
128 generate DNA fragment (single-hit² and double-hit³). The sequence bias estimates can be defined as the
129 ratio between the numbers of observed and expected DNase-seq reads starting at the middle of a particular
130 DNA sequence of length k (k -mer)¹⁵. We use here two approaches. The “DHS sequence bias” considers the
131 sequence bias estimates within DNase hypersensitive sites (DHSs) of each DNase-seq experiment. This
132 approach captures DNase I cleavage, read fragmentation and sequence complexity bias of DHSs of each
133 DNase-seq experiment^{7,15}. The “naked DNA sequence bias” considers the sequence bias estimates within
134 naked DNA DNase-seq experiments¹². In this case, all DNA regions are open, therefore the sequence bias
135 estimates will mainly capture the DNase I cleavage bias¹² (**Online Methods**). A clustering analysis of
136 sequence bias estimates forms two clear groups, which splits experiments from single-hit and double-hit
137 protocols (**Fig. 2**, **Supplementary Fig. 3**). This indicates that sequence biases are protocol-specific. Naked
138 DNA sequence bias estimates forms a sub-cluster within estimates from the double-hit experiments. This
139 highlights that DNase-seq experiments are influenced by more experimental artifacts than DNase sequence
140 cleavage bias alone.

141 Next, we extended the analysis by He et al.¹⁵ to evaluate the influence of sequence bias on all evaluated
142 footprinting methods based on the AUC at 10% false positive rate (FPR). HINT was evaluated with DNase-
143 seq signals corrected with either DHS sequence bias (HINT bias-corrected; HINT-BC) and naked DNA
144 sequence bias (HINT bias-corrected on naked DNase-seq; HINT-BCN). Our analysis shows that only six out

145 of 14 evaluated methods (Wellington, Neph, Boyle, DNase2TF, Centipede and FS-Rank) present a
146 significant negative Pearson correlation ($r = -0.35, -0.32, -0.28, -0.28, -0.24$ and -0.22 , respectively)
147 between their accuracy performance and amount of sequence bias (**Fig. 3a**; adjusted p -value < 0.05).
148 Equivalent results are also observed on the same TFs and cellular conditions analyzed in He et al.¹⁵
149 (**Supplementary Fig. 4**). Methods explicitly using 6-mer sequence bias statistics (HINT-BC, HINT-BCN and
150 FLR) or performing smoothing (Cuellar, BinDNase and PIQ) are not significantly influenced by sequence
151 bias. Moreover, the performance of HINT-BC is the least affected by sequence bias ($r = -0.06$). Pairwise
152 comparison of AUC at 10% FPR values of all three HINT variants (HINT-BC, HINT-BCN and HINT) indicates
153 significant gain in all predictions with sequence bias correction (adjusted p -value $< 10^{-30}$; **Supplementary**
154 **Fig. 5a**). There is no significant difference between HINT-BC and HINT-BCN, but we observe a higher AUC
155 on HINT-BC on all but seven TFs. This indicates an advantage of DHS sequence bias correction for the
156 footprint prediction problem.

157 As an example, we show sequence bias estimates, corrected and uncorrected DNase-seq average profiles
158 around TFBSs with highest AUC gain between HINT-BC and HINT (**Fig. 3b and c**; **Supplementary Fig. 6**).
159 The NRF1 and EGR1 DNase-seq profiles indicate that the bias-corrected signal fits better their sequence
160 affinity than the uncorrected signal. We observe that k -mers with high DHS sequence bias have a high CG
161 content ($r > 0.8$ in 11 out of 12 cell types; **Supplementary Fig. 7**). However, there is no significant
162 correlation between CG content of MPBSs and either AUC values or differences of AUC from HINT-BC,
163 HINT-BCN and HINT (p -value > 0.05 ; **Supplementary Fig. 5b**).

164 **Comparative analysis of footprinting methods**

165 Given its good performance^{10,15}, we evaluated the use of Tag Count (TC) as the ranking strategy instead of
166 each method's own ranking for BinDNase, Centipede, Cuellar, DNase2TF, FLR, PIQ and Wellington.
167 Previous to ranking by TC, site-centric methods required the definition of a minimum probability score to
168 define active footprints. In all cases, using TC yielded higher AUC values (10% FPR) than using their intrinsic
169 ranking metric (**Supplementary Fig. 8**). Concerning site-centric methods, the probability cutoff of 0.9 yielded
170 highest AUCs, with exception of BinDNase (best at 0.8). These parameters will be used in the next
171 evaluation analyses.

172 We next evaluated all the competing methods by measuring the AUC at 1%, 10% and 100% FPRs using the
173 TF ChIP-seq data. AUC at lower FPRs favors methods with higher sensitivity in expense of specificity. We
174 also estimated the AUPR, which is indicated for cases with imbalance of positive and negative examples²⁰,
175 and the FLR-Exp metric. Interestingly, all TF ChIP-seq based metrics indicate a very similar ranking ($r >$
176 0.98 ; **Fig. 4a**). There is also a high agreement between FLR-Exp and other metrics ($r > 0.88$). HINT-BC has
177 the highest FLR-Exp, AUC and AUPR values and significantly outperforms all methods with the exception of
178 HINT-BCN (adjusted p -value < 0.01 ; **Supplementary Fig. 9**; **Supplementary Tables 3-6**). Ignoring HINT
179 variants, the next top performing method is DNase2TF, which significantly outperforms all other methods
180 with the exception of PIQ (adjusted p -value < 0.01). PIQ outperforms all of its lower ranked competitors but
181 Wellington with AUC (1% FPR) and AUPR (adjusted p -value < 0.01). Concerning the performance of TC-
182 Rank, we observe that the AUC values for 10% and 100% FPR are very close to other footprinting methods
183 (**Fig. 4b**, **Supplementary Fig. 9**). This is not the case for AUC at 1% FPR or AUPR values. With the latter
184 statistics, all methods but Centipede and Cuellar have significant superior performance than TC (p -value $<$
185 0.01 ; **Supplementary Tables 3-6**).

186 **Transcription factor residence time**

187 Despite the high average prediction values of top performing footprint methods, they consistently perform
188 worst in a similar set of TFs, i.e. HINT-BC, DNase2TF and PIQ have 89% of TFs in common in the lower
189 quartile of AUC at 10% FPR (**Supplementary Dataset 1**). This list includes nuclear receptors, which has low
190 residence binding time⁷ and display a lower DNase I cleavage protection pattern (**Supplementary Fig. 10**).
191 To further investigate this, we propose a statistic inspired by the concepts presented in Sung et al.⁷ to detect
192 TFs with potential short residence time. The protection score measures the difference between the amounts

193 of DNase I digestion in the flanking regions and within the TFBS on bias-corrected DNase-seq signals. We
194 use this statistic to analyze the predictive performance of methods on TFs with distinct residence time. For
195 this, we used the comprehensive data set with 233 combinations of DNase-seq experiments and TFs (see
196 **Online Methods**).

197 We observed that TFs with known short residence time on DNA, such as nuclear receptors AR²¹, ER²² and
198 GR²³, present a negative protection score (**Fig. 5a**). TFs with intermediate and long residence time on DNA
199 (C-JUN²⁴ and CTCF²⁵, respectively) present a positive protection score. The amount of protection is clearly
200 reflected in the bias-corrected DNase-seq profiles (**Fig. 5b-d**). In addition, **Figure 5a** also reveals an
201 association of the protection score and the AUC of HINT-BC. Overall, the protection score positively
202 correlates with the AUC values of evaluated methods, such as TC ($r = 0.19$) and HINT-BC ($r = 0.26$), and
203 negatively correlates ($r = -0.49$) with the sequence bias (adjusted p -value < 0.05). These results reinforce
204 the concept that TFs with potential short residence time can be poorly detected via DNase-seq footprints.

205 **DISCUSSION**

206 Our comparative evaluation analysis indicates the superior performance (in decreasing order) of HINT,
207 DNase2TF and PIQ in the prediction of active TFBSs in all evaluated scenarios. Moreover, tools
208 implementing these methods were user friendly and had lower computational demands than other evaluated
209 methods. Clearly, the choice of computational footprinting approaches should also be based on experimental
210 design aspects. For example, PIQ is the only method supporting analysis of replicates and time-series. On
211 the other hand, studies requiring footprint predictions for latter *de novo* motif analysis should use
212 segmentation approaches as HINT or DNase2TF. In contrast to positive evaluations of the TC-Rank by
213 previous works^{10,15}, we show that it has poor sensitivity performance as indicated by the AUC at low FPR
214 levels. On the other hand, the TC statistic provides the best strategy to rank footprint predictions from other
215 methods.

216 The refined DNase-seq protocol and experimental artifacts presented in He et al.¹³ and TF binding time
217 presented in Sung et al.⁷ underscore that robust *in silico* techniques are required to correct for experimental
218 artifacts and to derive valid biological predictions. The correction of DNase-seq signal with DHS sequence
219 bias estimates virtually removes the effects of sequence bias artifacts on computational footprinting. We
220 demonstrated that such correction can be performed prior to the execution of the computational footprinting
221 method. On the other hand, ignoring experimental artifacts might lead to false predictions, as observed
222 previously for predicted *de novo* motifs (**Supplementary Fig. 11**). Moreover, the simple protection score can
223 indicate footprints of TFs with potential short binding time. Thus, footprint predictions of TFs with low
224 protection score should be interpreted with caution.

225 The assessment of footprint methods is a demanding task, both computationally and technically. We have
226 created a fair and reproducible benchmarking data set for evaluation of protein-DNA binding using two
227 validation approaches: TF ChIP-seq and FLR-Exp. Although the rationales of the ChIP-seq and FLR-Exp
228 evaluation procedures are, in principle, very different, we observed a high agreement between their
229 respective ranking of methods. This is evidence that this study provides a robust map of the accuracy of
230 state-of-the-art computational footprinting methods. Finally, this study provides all statistics, basic data and
231 scripts to evaluate future computational footprinting methods. This is an important resource for increasing
232 transparency and reproducibility of research on computational methods for DNase-seq data.

233 **Acknowledgements:** This work was supported by the Interdisciplinary Center for Clinical Research (IZKF
234 Aachen), RWTH Aachen University Medical School, Aachen, Germany (to E.G.G., M.A. and I.G.C.) and the
235 Excellence Initiative of the German Federal and State Governments and the German Research Foundation
236 through Grant GSC 111 (M.A. and I.G.C.).

237
238 **Author Contributions:** E.G.G., M.Z. and I.G.C. designed the research. E.G.G. wrote HINT program code.
239 E.G.G., M.A. and I.G.C. analyzed data. E.G.G., M.Z. and I.G.C. wrote the manuscript.

240
241 **Competing Financial Interests:** The authors declare no competing financial interests.

- 242 ¹ The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- 243 ² Crawford, G.E. et al. Genome-wide mapping of DNase hypersensitive sites using massively parallel signature sequencing (MPSS). *Genome Res.* **16**, 123–131 (2006).
- 244 ³ Sabo, P.J. et al. Genome-wide identification of DNase I hypersensitive sites using active chromatin sequence libraries. *Proc. Natl. Acad. Sci. USA* **101**, 4537–4542 (2004).
- 245 ⁴ Neph, S. et al. An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489**, 83–90 (2012).
- 246 ⁵ Boyle, A.P. et al. High-resolution genome-wide in vivo footprinting of diverse transcription factors in human cells. *Genome Res.* **21**, 456–464 (2011).
- 247 ⁶ Piper, J. et al. Wellington: a novel method for the accurate identification of digital genomic footprints from DNase-seq data. *Nucleic Acids Res.* **41**, e201 (2013).
- 248 ⁷ Sung, M.-H.H., Guertin, M.J., Baek, S. & Hager, G.L. DNase footprint signatures are dictated by factor dynamics and DNA sequence. *Mol. Cell* **56**, 275–285 (2014).
- 249 ⁸ Gusmao, E.G., Dieterich, C., Zenke, M. & Costa I.G. Detection of active transcription factor binding sites with the combination of DNase hypersensitivity and histone modifications. *Bioinformatics* **30**, 3143–3151 (2014).
- 250 ⁹ Pique-Regi, R. et al. Accurate inference of transcription factor binding from DNA sequence and chromatin accessibility data. *Genome Res.* **21**, 447–455 (2011).
- 251 ¹⁰ Cuellar-Partida, G. et al. Epigenetic priors for identifying active transcription factor binding sites. *Bioinformatics* **28**, 56–62 (2012).
- 252 ¹¹ Sherwood, R.I. et al. Discovery of directional and nondirectional pioneer transcription factors by modeling DNase profile magnitude and shape. *Nat. Biotechnol.* **32**, 171–178 (2014).
- 253 ¹² Yardimci, G.G., Frank, C.L., Crawford, G.E. & Ohler, W. Explicit DNase sequence bias modeling enables high-resolution transcription factor footprint detection. *Nucleic Acids Res.* **42**, 11865–11878 (2014).
- 254 ¹³ Kähärä, J. & Lähdesmäki, H. BinDNase: A discriminatory approach for transcription factor binding prediction using DNase I hypersensitivity data. *Bioinformatics* **31**, 2852–2859 (2015).
- 255 ¹⁴ Stergachis, A.B. et al. Conservation of trans-acting circuitry during mammalian regulatory evolution. *Nature* **515**, 365–370 (2014).
- 256 ¹⁵ He, H.H. et al. Refined DNase-seq protocol and data analysis reveals intrinsic bias in transcription factor footprint identification. *Nat. Methods* **11**, 73–78 (2014).
- 257 ¹⁶ Meyer, C. & Liu, X. Identifying and mitigating bias in next-generation sequencing methods for chromatin biology. *Nat. Rev. Genet.* **15**, 709–721 (2014).
- 258 ¹⁷ Park, P.J. ChIP-seq: advantages and challenges of a maturing technology. *Nat. Rev. Genet.* **10**, 669–680 (2009).
- 259 ¹⁸ Teytelman, L., Thurtle, D.M., Rine, J. & Oudenaarden, A. Highly expressed loci are vulnerable to misleading ChIP localization of multiple unrelated proteins. *Proc. Natl. Acad. Sci. USA* **110**, 18602–18607 (2013).
- 260 ¹⁹ Editorial. The difficulty of a fair comparison. *Nat. Methods* **12**, 273–273 (2015).
- 261 ²⁰ Davis, J. & Goadrich, M. The relationship between Precision-Recall and ROC curves. *Proceedings of the 23rd international conference on Machine learning – ICML 2006*, 233–240 (2006).
- 262 ²¹ Tewari, A.K. et al. Chromatin accessibility reveals insights into androgen receptor activation and transcriptional specificity. *Genome Biol.* **13**, R88 (2012).
- 263 ²² Sharp, G.D. et al. Estrogen-receptor-alpha exchange and chromatin dynamics are ligand- and domain-dependent. *J. Cell Sci.* **119**, 4101–4116 (2006).
- 264 ²³ McNally, J.G., Müller, W.G., Walker, D., Wolford, R. & Hager, G.L. The glucocorticoid receptor: rapid exchange with regulatory sites in living cells. *Science* **287**, 1262–1265 (2000).
- 265 ²⁴ Malnou, C.E. et al. Heterodimerization with different Jun proteins controls c-Fos intranuclear dynamics and distribution. *J. Biol. Chem.* **285**, 6552–6562 (2010).
- 266 ²⁵ Nakahashi, H. et al. A genome-wide map of CTCF multivalency redefines the CTCF code. *Cell Rep.* **3**, 1678–1689 (2013).

286 **Figure 1 | FLR-Exp evaluation metric.** (a) FLR score distribution of footprints predicted with HINT-BC overlapping with MPBSs of selected TFs. These TFs have increasing expression in K562 (red) compared with H1-hESC cells (blue). The signed Kolmogorov-Smirnov (KS) statistic quantifies the separation of both distributions. The box plot depicts the distribution median value (middle dot) and first and third quartiles (box extremities). The whiskers represent the 1.5 IQR and external dots represent outliers (data greater than or smaller than 1.5 IQR). (b) Scatter plot with signed KS statistic and expression fold change for 143 TFs. There is a clear association between TF expression and KS statistic ($r = 0.97$, adjusted p -value $< 10^{-10}$).

294 **Figure 2 | Clustering of bias estimates.** Ward's minimum variance clustering based on pairwise Spearman correlation coefficient (r) from bias estimates of all ENCODE's Tier 1 and naked DNA DNase-seq data. DNase-seq experiments were based on single-hit (red), double hit (blue) protocols or naked DNA (yellow).

298 **Figure 3 | Effects of sequence biases on methods.** (a) Association between the performance of footprinting methods (relative to TC-Rank performance) and their sequence bias estimated for 88 TFs binding on cell types H1-hESC and K562. The x-axis represents the correlation between the uncorrected and bias signal (observed versus bias signal; OBS). The OBS is evaluated for each TF by measuring the uncorrected DNase-seq signal and the bias signal for every MPBS that overlaps a footprint from the evaluated method. Then, the Spearman correlation is evaluated between the average uncorrected and bias signals. Higher OBS values indicate higher bias. The y-axis represents the ratio between the AUC at 10% FPR for each evaluated method and the TC-Rank method; higher values indicate higher accuracy. (b-c) Average bias signal (top) and uncorrected/bias-corrected DNase-seq signal (bottom) for the TFs: (b) NRF1 and (c) EGR1. Signals in the bottom graph were standardized to be in the interval [0,1]. The motif logo represents all underlying DNA sequences centered on the TFBSs.

310 **Figure 4 | Evaluation of computational footprinting methods.** (a) Average rankings for the evaluated
 311 computational footprinting methods. The rankings are given for all evaluation criteria: FLR-Exp, TF ChIP-seq
 312 based AUC (at 100%, 10% and 1% FPR) and AUPR. (b) For all evaluated methods we show the FLR-Exp
 313 values (as a combination of all pairwise comparison within cell types H1-hESC, K562 and GM12878),
 314 median TF ChIP-seq based AUC (at 100%, 10% and 1% FPR) values and median AUPR values. HINT-BC,
 315 HINT-BCN, HINT, DNase2TF are ranked as top four methods by all evaluation metrics. All baseline methods
 316 (FS-Rank, PWM-Rank and TC-Rank) are in the bottom four positions of the ranks. Note that BinDNase could
 317 not be evaluated with the FLR-Exp, as it requires ChIP-seq data for training.
 318

319 **Figure 5 | Impact of transcription factor residence binding time on computational footprinting.** (a)
 320 Scatter plot with the protection score (*x*-axis) versus TF ChIP-seq based AUC (at 10% FPR) of HINT-BC (*y*-
 321 axis) for 233 TFs binding on 11 cell types. We highlight nuclear receptors AR, ER and GR (short residence
 322 time, red); C-JUN (intermediate residence time, blue); CTCF (long residence time, green) and TFs with either
 323 high (> 6) protection score or low (< 0.8) AUC values (grey). (b-d) Average bias signal (top) and
 324 uncorrected/bias-corrected DNase-seq signal (bottom) for the TFs (b) ER, (c) C-JUN and (d) CTCF. Signals
 325 in the bottom graph were standardized to be in the interval [0,1]. The motif logo represents all underlying
 326 DNA sequences centered on the TFBSs.
 327

328 **Table 1 | Overview of methods.** Main characteristics of the evaluated methods. Methods obtain a '+' sign
 329 for availability if they are public available. Boyle method is not public, but authors provide footprint
 330 predictions of a few cells. Concerning usability, methods natively supporting standard genomic files and
 331 being executed with few commands (≤ 3) display a '+' sign.
 332

Name	Type	Algorithm	Bias Correction	Resolution/Smoothing	Footprint Ranking	Availability	Usability	Others
BinDNase	SC	Logistic regression	None	Base pair / sliding window	Probability	+	-	Require TF ChIP-seq for training
Boyle	SEG	HMM	None	Base pair	None	-	-	
Centipede	SC	Bayesian mixture model	None	Base pair	Probability	+	-	Integrates histone and sequence data
Cuellar	SC	Weighted motif match	None	Sliding window	PWM score	+	-	
DNase2TF	SEG	Sliding window	4-mer (DHS sequence bias)	Base pair	<i>p</i> -values	+	+	
FLR	SC	Mixture model	6-mer (naked DNA sequence bias)	Base pair	Log-odds	+	-	Bias correction for each TF
HINT	SEG	HMM	6-mer (DHS sequence bias)	Base pair	TC	+	+	Integrates histones

Neph	SEG	Sliding window	None	Base pair	FS	+	-	
PIQ	SEG	GP / expectation propagation	None	Base pair / GP	Probability	+	+	Support replicates, time series
Wellington	SEG	Sliding window	None	Base pair	p-value	+	+	

333

334

METHODS

335

336

337

338

339

340

341

342

343

344

Data. DNase-seq aligned reads were obtained from ENCODE¹. To perform the computational footprint experiments, we obtained data regarding cell types H1-hESC, HeLa-S3, HepG2, Huvec, K562, LNCaP and MCF-7 from Crawford's Lab (labeled with the initials of their institution "DU") and cell types H7-hESC, HepG2, Huvec, K562 and m3134 from Stamatoyannopoulos' lab (labeled with the initials of their institution "UW"). We also used naked DNA (deproteinized) DNase-seq experiments from cell types MCF-7 and K562 (DU)¹² and IMR90 (UW)²⁶. DNase-seq experiments labeled with "DU" follow the single-hit protocol, while the experiments labeled with "UW" follow the double-hit protocol. In addition, to perform the DNase-seq bias estimation clustering, we used all cell types from ENCODE's Tier 1 and Tier 2 cell types¹. See **Supplementary Table 2** for a full DNase-seq data description.

345

346

347

348

349

350

351

352

353

354

355

356

Transcription factor (TF) ChIP-seq enriched regions (peaks and summits) were obtained in ENCODE analysis working group (AWG)¹ track with exception of the following experiments, in which the enriched regions were obtained using bowtie-2²⁷ and MACS²⁸. AR (R1881 treatment) ChIP-seq raw sequences for LNCaP cell type was obtained in gene expression omnibus (GEO) with accession number GSM353644²⁹. ER (40 and 160 minutes after estradiol treatment) ChIP-seq raw sequences for MCF-7 cell type was obtained in GEO with accession number GSE54855³⁰. GR (dexamethasone treatment) ChIP-seq raw sequences for m3134 cell type was obtained in the sequence read archive (SRA) under study number SRP004871³¹. All organism-specific data (DNase-seq and ChIP-seq) are based on the human genome build 37 (hg19), except the DNase-seq for m3134 and ChIP-seq for GR, which were based on mouse genome build 37 (mm9). Chromosome Y was removed from all analyses. Expression of cells H1-hESC, K562 and GM12878 were obtained from ENCODE (GSE12760 and GSE14863)¹.

357

358

359

360

361

362

363

TF motifs (position frequency matrices; PFMs) were obtained from the Jaspar³², Uniprobe³³ and Transfac³⁴ repositories. Non-organism-specific data (PFMs) were obtained for the subphylum Vertebrata. *De novo* PFMs 0458 and 0500 were downloaded from ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/supplementary/integration_data_jan2011/byDataType/footprints/jan2011/de.novo.pwm⁴. The accession codes for all TF ChIP-seq experiments and PFM IDs are available in the **Supplementary Datasets 1a and 2b-d**.

364

365

366

367

368

369

Sequence bias correction. *DNase I hypersensitivity sites.* A first task is the identification of DNase I hypersensitivity sites (DHSs). A nucleotide-resolution genome-wide signal was created for each DNase-seq data set by counting reads mapped to the genome. Here, we considered only the 5' position of the aligned reads (position at which DNase I cleaved the DNA). The genomic signal was created by counting the number of reads that overlapped at each genomic position.

370

371

More formally, we define a raw genomic signal as a vector

$$x = \langle x_1, \dots, x_N \rangle,$$

372

373

374

375

where N equals the number of bases in the genome and each $x_i \in N^0$ is the number of DNase-seq reads in which the 5' position mapped to position i . We also generate strand specific counts X^s , where $s \in \{+, -\}$ describes the strand the read was mapped to.

376

377

378

379

DHSs are estimated based on the DNase I raw signal. First, the F-seq software³⁵ was used to create smoothed DNase-seq signals using Parzen density estimates. Then, the smoothed signal x^{fseq} was fit to a gamma distribution,

$$x^{fseq} \sim \Gamma(\kappa, \theta),$$

380 by evaluating κ and θ based on mean and standard deviation estimates. Finally, the enriched regions
 381 (DHSs) were found by establishing a cutoff based on a p -value of $0.01^{1,35}$. We refer to DHSs as a set of
 382 genomic intervals

$$383 \quad H = \{h_1, \dots, h_L\},$$

384 where $h_i = [m, n]$ for $m < n \in N$ and L is the total number of DHSs. We ignore for simplicity of notation the
 385 fact that intervals are defined on distinct chromosomes or contigs.

386
 387 *Estimation of DNase-seq sequence bias.* We use two approaches to estimate sequence bias of DNase-seq
 388 experiments: (1) aligned reads inside DHSs from DNase-seq experiments (termed ‘‘DHS sequence bias’’)
 389 following He et al.¹⁵ and (2) all aligned reads for naked DNA experiments (termed ‘‘naked DNA sequence
 390 bias’’) following Yardımcı et al.¹². The observed cleavage score for a k -mer w corresponds to the number of
 391 DNase I cleavage sites centered at w . The background cleavage score is defined by the total number of
 392 times w occurs. Then, the bias estimation is computed as the ratio between the observed and background
 393 cleavage scores. Mathematical formalizations of the bias estimation will be made based on the DHS
 394 sequence bias approach.

395
 396 We define G^s as the reference genome sequence with length N for strand $s \in \{+, -\}$. $G^s[i..j]$ indicates the
 397 sequence from positions i to j (including both within the interval). For each k -mer w with length k the
 398 observed cleavage score o_w can be calculated as

$$399 \quad o_w^s = 1 + \sum_{i=1}^L \sum_{j \in h_i} x_j^s \mathbf{1} \left(G^s \left[j - \frac{k}{2} .. j + \frac{k}{2} \right] = w \right),$$

400 where $\mathbf{1}(\cdot)$ is an indicator function.

401

402 Similarly, the background cleavage score r_w can be evaluated as

$$403 \quad r_w^s = 1 + \sum_{i=1}^L \sum_{j \in h_i} \mathbf{1} \left(G^s \left[j - \frac{k}{2} .. j + \frac{k}{2} \right] = w \right).$$

404

405 Finally, the cleavage bias b_i^s for a genomic position $k + 1 \leq i \leq N - k + 1$, given that $w = G^s \left[i - \frac{k}{2} .. i + \frac{k}{2} \right]$, can be calculated as

$$407 \quad b_i^s = \frac{o_w^s \cdot R}{r_w^s \cdot O^s},$$

408 where O^s indicates the total number of reads aligned to strand s in DHSs

$$409 \quad O^s = \sum_{i=1}^L \sum_{j \in h_i} x_j^s,$$

410 and R indicates the total number of k -mers in DHS positions

$$411 \quad R = \sum_{i=1}^L \sum_{j \in h_i} 1.$$

412

413 The bias score b_i^s represents how many times the k -mer sequence $G^s \left[i - \frac{k}{2} .. i + \frac{k}{2} + 1 \right]$ was cleaved by the
 414 DNase I enzyme in comparison to its total occurrence in: (1) DHSs (DHS sequence bias approach); (2) the
 415 entire genome (naked DNA sequence bias approach). As observed by He et al.¹⁵ a 6-mer bias model
 416 captures more information than $k < 6$ models and the information added with $k > 6$ models are not
 417 significant. Therefore, in this study, all analyses were performed using a 6-mer bias model.

418

419 *DNase-seq sequence bias correction.* A ‘‘smoothed corrected signal’’ was calculated using smoothed
 420 versions of both raw DNase-seq (\hat{x}_i^s) and the bias score signal (\hat{b}_i^s)¹⁵. These smoothed signals were based
 421 on a 50 bp window and can be written as

$$422 \quad \hat{x}_i^s = \sum_{j=i-25}^{i+24} x_j^s$$

$$423 \quad \hat{b}_i^s = \frac{b_i^s}{\sum_{j=i-25}^{i+24} b_j^s}.$$

423

424 With these results we are able to define the smoothed corrected signal as

$$425 \quad c_i^s = \hat{x}_i^s \hat{b}_i^s.$$

426

427 Finally, the bias-corrected DNase-seq genomic signal (y) can be obtained by applying

$$428 \quad y_i^s = \log(x_i^s + 1) - \log(c_i^s + 1). \quad (1)$$

429

430 The corrected DNase-seq signal generated by **equation (1)** may include negative values. Since some
431 posterior statistical analyses required a signal consisting only of positive values, we have shifted the entire
432 signal by adding the global minimum value.

433

434 **Computational footprinting methods.** In this section we present an overview of the computational
435 footprinting methods used in this study. Computational resources necessary to the execution of each method
436 were summarized in **Supplementary Table 1**.

437

438 *Neph method.* Neph et al.⁴ used a simplified version of the segmentation method originally proposed in
439 Hesselberth et al.³⁶. Their method consists on applying a sliding window to find genomic regions (6-40 bp)
440 with low DNase I cleavage activity between regions (3-10 bp) with intense DNase I digestion. The footprint
441 score (FS) is evaluated and used to determine the most significant predictions.

442

443 We obtained the footprint predictions for cell type K562 (DU) in
444 ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/supplementary/integration_data_jan2011/byDataType/footprints/jan2011/all.footprints.gz⁴. As predictions were not available for other DNase-seq experiments, we
445 obtained the scripts and parameterization through Neph et al.⁴ footprinting method code repository at
446 <https://github.com/StamLab/footprinting2012>. Briefly, we used the DNase I raw signal as input with the
447 parameters from the original publication: flanking component length varied between 3-10 bp and central
448 footprint region length varied between 6-40 bp. Afterwards, the footprints were filtered by an FDR of 1%,
449 which was estimated based on the FS distribution in each cell type⁴. Finally, we consider only predictions
450 that occurred within DNase-seq hotspots, evaluated using the method first described in Sabo et al.³⁷. We
451 obtained all hotspots generated by Stamatoyannopoulos' lab in ENCODE¹ for cell types GM12878
452 (wgEncodeEH000492; GSM736496 and GSM736620), H1-hESC (wgEncodeEH000496; GSM736582) and
453 K562 (wgEncodeEH000484; GSM736629 and GSM736566). We will refer to this framework as “Neph”.

454

455 *Boyle method.* Boyle et al.⁵ designed a segmentation approach, which is based on using hidden Markov
456 models (HMMs) to predict footprints in specific DNase I cleavage patterns. Briefly, the HMM uses a
457 normalized DNase-seq cleavage signal to find regions with depleted DNase I digestion (footprints) between
458 two peaks of intense DNase I cleavage. Such pattern reflects the inability of the DNase I nuclease to cleave
459 sites where there are proteins bound. As the DNase-seq profiles required a nucleotide-resolution signal,
460 which is usually noisy, the authors used a Savitzky-Golay smoothing filter to reduce noise and to estimate
461 the slope of the DNase-seq signal³⁸. Their HMM had five states, with specific states to identify the
462 decrease/increase of DHS signals around the peak-dip-peak region. Since no source code or software is
463 provided, we used footprint predictions from Boyle et al.⁵ available at
464 <http://fureylab.web.unc.edu/datasets/footprints/>. We will refer to this method as “Boyle”.

465

466 *Centipede.* Centipede is a site-centric approach, which gathers experimental and genomic information
467 around motif-predicted binding sites (MPBSs). It then uses a Bayesian mixture model approach to label each
468 retrieved site as 'bound' or 'unbound'⁹. The experimental and genomic data used include DNase-seq,
469 position weight matrix (PWM) bit-score, sequence conservation and distance to the nearest transcription
470 start site (TSS). The experimental data input was generated by fetching the raw DNase-seq signal
471 surrounding a 200 bp window centered on each MPBS. Additionally, to create the genomic data input, we
472 obtained PhastCons conservation score (placental mammals on the 46-way multiple alignment)³⁹ and
473 Ensembl gene annotation from ENCODE^{1,40} to create the prior probabilities in addition to the PWM bit-score.

474

475 Centipede software was obtained at <http://centipede.uchicago.edu/> and executed to generate posterior
476 probabilities of regions being bound by TFs. We have previously observed that Centipede is sensitive to
477 certain parameters. Therefore, Centipede parameterization was defined with an extensive computational
478 evaluation described in Gusmao et al.⁸.

479

480 *Cuellar Method.* Cuellar-Partida et al.¹⁰ proposed a site-centric method to include DNase-seq data as priors
481 for the detection of active transcription factor binding sites (TFBSs). It is based on a probabilistic
482 classification approach to compute better log-posterior odds score than the ones observed by purely
483 sequence-based approaches. We applied this method as described in Cuellar-Partida et al.¹⁰. We created a
484 smoothed DNase-seq input signal by evaluating the number of DNase-seq cleavage based on a 150 bp
485 window with 20 bp steps. We obtained their scripts at <http://research.imb.uq.edu.au/t.bailey/SD/Cuellar2011/>
486 and created priors using the smoothed version of the DNase-seq signal. As suggested by the authors, the
487 priors were submitted to the program FIMO⁴¹ to obtain the predictions. We will refer to this method as
488 “Cuellar”.

489

490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549

Wellington. Wellington is a segmentation approach based on a Binomial test. For a given candidate footprint, it tests the hypothesis that there are more reads in the flanking regions than within the footprint. Following an observation that DNase-seq cuts of the double-hit protocol are strand-specific, Wellington only considers reads mapped to the upstream flanking region of the footprints. Wellington automatically detects the size of footprints (within a user-defined interval) and sets flanking regions at a user-defined length. We have obtained Wellington's source code in <http://jpiper.github.com/pyDNase> and executed it with default parameters. Briefly, we used a footprint FDR cutoff of -30 , footprint sizes varying between 6 and 40 with 1 bp steps and shoulder size (flanking regions) of 35 bp.

Protein interaction quantification (PIQ). The protein interaction quantification (PIQ) is a site-centric method, which uses Gaussian process to model and smooth the footprint profiles around candidate MPBSs (± 100 bp)¹¹. Active footprints are estimated with an expectation propagation algorithm. Finally, PIQ indicates the set of motifs which footprint signals are distinguishable from noise to reduce the set of candidate TFs. We obtained PIQ implementation in <http://piq.csail.mit.edu> and executed it with default parameters, which can be found in the script *common.r*. Briefly, MPBSs were generated with the script *pwmmatch.exact.r*. The DNase-seq signal was created using the script *bam2rdata.r*. And the footprints were detected with the script *perft.r*.

Footprint mixture (FLR). Yardımcı et al.¹² proposed a site-centric method based on a mixture of multinomial models to detect active/inactive MPBSs. The method uses an expectation maximization algorithm to find a mixture of two multinomial distributions, representing active (footprints) and inactive (background) MPBSs. The background model is initialized with either naked DNA sequence bias frequencies or estimated *de novo*. After successful estimation, MPBSs are scored with the log odds ratio for the footprint *versus* background model. The model takes DNase-seq cuts within a small window around the candidate profiles (± 25 bp) as input. DNase-seq sequence bias is estimated for 6-mers based on the DNA sequences extracted within the same regions in which the cuts were retrieved. Method implementation was obtained in https://ohlerlab.mdc-berlin.de/software/FootprintMixture_109/. We executed the method using naked DNA sequence bias frequencies for initialization of the background models. The width of the window surrounding the TFBS (*PadLen*) was set to the default value of 25 bp. Also, we use the expectation maximization to re-estimate background during training (argument *Fixed* set to *FALSE*). We will refer to this method as "FLR".

DNase2TF. DNase2TF is a segmentation approach based on a binomial z-score, which evaluates the depletion of DNase-seq reads around the candidate footprints⁷. At a second step, DNase2TF interactively merges close candidate footprints whenever they improve depletion scores. DNase2TF corrects for DNase I sequence bias using cleavage statistics for 2- or 4-mers. We obtained source code from <http://sourceforge.net/projects/dnase2tfr/> and executed DNase2TF with a 4-mer sequence bias correction. Other parameters were set to their default values: *minw* = 6, *maxw* = 30, *z_threshold* = -2 and *FDR* = 10^{-3} .

HINT, HINT-BC and HINT-BCN. Recently, Gusmao et al.⁸ have proposed the segmentation method HINT (HMM-based identification of transcription factor footprints) as an extension of Boyle method⁵. HINT is based on eight-state multivariate HMMs and combines DNase-seq and histone modification ChIP-seq profiles at the nucleotide level for the identification of footprints. The pipeline of HINT method starts by normalizing the DNase I cleavage signal using within- and between-dataset normalizations. Then, the slope of the normalized signals is evaluated to identify the DNase-seq signal increase and decrease. Afterwards, an HMM is trained on a supervised manner (maximum likelihood) based on a single manually annotated genomic region. To aid such manual annotation the normalized and slope signals are used in combination with MPBSs for all available PFMs in the repositories Jaspar³² and Uniprobe³³. Finally, the Viterbi algorithm is performed on the trained HMMs inside regions consisting of DHSs extended by 5,000 bp upstream and downstream. All parameters were set as described in Gusmao et al.⁸.

We have performed two modifications to the method described in Gusmao et al.⁸. First, to perform a standardized comparison, we modified HINT to allow only DNase-seq data. The modified HMM model contains five states. The three histone-level states were removed and new transitions were created from the *BACKGROUND* state to the *DNase UP* state and from the *DNase DOWN* state to the *BACKGROUND* state. The second modification concerns the use of bias-corrected DNase-seq signal prior to normalization steps. We will call the method HINT bias-corrected (HINT-BC), for correction based on "DHS sequence bias", and HINT bias-corrected on naked DNA (HINT-BCN) for the "naked DNA sequence bias" estimation. These modifications required retraining of the HMM models. For this, we used the same manual annotation described in Gusmao et al.⁸. The novel methods and trained models are available as a command-line tool at www.costalab.org/hint-bc.

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603

BinDNase. BinDNase is a site-centric method based on logistic regression to predict active/inactive MBBSs¹³. The algorithm starts with base pair resolution DNase-seq signal around the MPBSs (± 100 bps) and selects discriminatory features using a backward greedy approach. As a supervised approach, the method requires positive and negative examples, which can be obtained from TF ChIP-seq data. We have used DNase-seq data around MPBSs on chromosome 1 for training. These MPBSs were subsequently removed from the evaluation procedure. The definition of positive and negative examples was the same as in our evaluation data sets. Note that this is the only method evaluated here which requires TF ChIP-seq examples for training. We also point the fact that BinDNase did not successfully executed for 19 TFs of our evaluation data set (POU5F1, REST, RFX5, SP1, SP2, SRF, TCF12 and ZNF143 binding in H1-hESC; ARID3A, CTCF, IRF1, MEF2A, PU1, REST, RFX5, SP1, SP2, STAT2 and ZNF263 binding in K562) given our maximum running time criteria (three weeks). Method implementation was obtained at <http://research.ics.aalto.fi/csb/software/bindnase/> and required/provided no parameter selection.

Footprint score rank (FS-Rank). He et al.¹⁵ used a site-centric MPBS ranking scheme termed “footprint score (FS)”, which is based on a scoring metric from the footprinting methodology proposed in Neph et al.⁴. The FS statistic is defined as

$$FS_{MPBS_i} = - \left(\frac{n_{C,i+1}}{n_{R,i+1}} + \frac{n_{C,i+1}}{n_{L,i+1}} \right),$$

where $MPBS_i = [m_i, n_i]$ is the i -th MPBS which extends from genomic positions m_i to n_i and $\overline{MPBS}_i = (m + n)/2$. The FS uses the DNase-seq signal in the center ($n_{C,i}$) of the MPBS and its upstream ($n_{L,i}$) and downstream ($n_{R,i}$) flanking regions. These variables can be defined as

$$\begin{aligned} n_{C,i} &= \sum_{j=m_i}^{n_i} x_j \\ n_{R,i} &= \sum_{j=n_i}^{2n_i-m_i} x_j \\ n_{L,i} &= \sum_{j=2m_i-n_i}^{m_i} x_j \end{aligned} \quad (2)$$

Tag count rank (TC-Rank). The site-centric method which we refer to as “tag count (TC)”, corresponds to the number of DNase I cleavage hits in a 200 bp window around predicted TFBS as defined in He et al.¹⁵. This can be written as

$$TC_{MPBS_i} = \sum_{j=\overline{MPBS}_i-100}^{\overline{MPBS}_i+99} x_j.$$

Both TC and FS can be used as quality scores for footprints. However, as a method (termed TC-Rank and FS-Rank) it consists on attributing these quality scores to each MPBS and evaluating the performance at these ranked MPBS. This observation also holds for the PWM-Rank method described below.

Evaluation. Motif-predicted binding sites (MPBSs). Method evaluation was performed with a site-centric binding site statistics. For this, we generated position weight matrices (PWMs) from PFMs by evaluating the information content of each position and performing background nucleotide frequency correction⁴². This was performed using Biopython⁴³. Then, we created MPBSs by matching all PWMs against the human (hg19) and mouse (mm9) genomes using the fast performance motif matching tool MOODS⁴⁴. This procedure produces “PWM bit-scores” for every match. We determined a bit-score cutoff threshold by applying the dynamic programming approach described in Wilczynski et al.⁴⁵ with a false positive rate (FPR) of 10^{-4} . All site-centric scores were based on the set of MPBSs after the application of the cutoff threshold. Also, the PWM bit-score was used as a baseline method and will be referenced as “PWM-Rank”.

Method comparison. Methods were evaluated using a site-centric approach¹⁰, which combines MPBSs with ChIP-seq data for every TF. In this scheme, MPBSs with ChIP-seq evidence (located within 100 bp from the ChIP-seq peak summit) are considered “true” TFBSs; while MPBSs without ChIP-seq evidence are considered “false” TFBSs. Every TF prediction that overlaps a true TFBS is considered a correct prediction (true positive; TP) and every prediction that overlaps with a false TFBS is considered an incorrect prediction (false positive; FP). Therefore, true negatives (TN) and false negatives (FN) are, respectively, false and true TFBSs without overlapping predictions. To assess the accuracy of digital genomic footprinting methods we created receiver operating characteristic (ROC) curves. Briefly, ROC curves describe the sensitivity (recall) increase as we decrease the specificity of the method. The area under the ROC curve (AUC) metric was evaluated at 100%, 10% and 1% false positive rates (FPRs). We also evaluated the area under the precision-recall curve (AUPR). This metric is indicated for problems with imbalanced data sets (distinct number of positive and negative examples)^{20,46}.

604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662

Segmentation approaches (Boyle, DNase2TF, HINT, Neph and Wellington) provide footprint predictions that do not necessarily encompass all MPBSs. To create full ROC curves for these methods, we first ranked all predicted sites by their DNase I cleavage tag count followed by all non-predicted sites ranked by their tag count. In order to present a fair comparison, this approach was also applied to all site-centric methods (Centipede, Cuellar, FLR and PIQ). For that, we considered distinct probability thresholds of (0.8, 0.85, 0.9, 0.95, 0.99) for detection of footprints on all site-centric methods. We performed additional experiments to select the best threshold per method (see **Supplementary Fig. 8**).

Our TF ChIP-seq based comparative experiments comprise the following three evaluation scenarios. All evaluation statistics and method performances are available at the **Supplementary Dataset 1**.

He Dataset: To replicate the analysis performed by He et al.¹⁵, we analyzed DNase-seq from cell types K562 (UW), LNCaP (DU) and m3134 (UW) on 36 TFs and we evaluated the methods PWM, FS, TC, HINT, HINT-BC and HINT-BCN.

Benchmarking Dataset: For comparative analysis of several competing methods, we selected the two cell types with highest number of TF ChIP-seq data sets evaluated in our study: K562 (DU) with 59 TFs and H1hesc (DU) with 29 TFs. We can therefore make use of predictions provided by Gusmao et al.⁸ and Boyle et al.⁵, which includes evaluation of PWM, Boyle, Cuellar, Centipede, HINT and Neph methods. For this data set, we have estimated novel footprints for FS, TC, HINT-BC, HINT-BNC, DNase2TF, PIQ, Wellington and FLR methods, which were not previously evaluated.

Comprehensive dataset: Lastly, we have compiled a comprehensive data set containing 233 combinations of cells and TFs with matching cellular background. This data set was built from a catalog of 144 TF ChIP-seq and 13 DNase-seq data sets. This data is used to evaluate the effects of bias correction and TF binding time. In this scenario we evaluated the methods PWM, FS, TC, HINT, HINT-BC and HINT-BCN.

Expression-based evaluation (FLR-Exp). As shown in Yardımcı et al.¹², ChIP-seq evaluation of putative TFBSs may present biases regarding the fact that ChIP-seq data alone is not able to distinguish direct from indirect binding events. Consequently, we performed an evaluation procedure which combines MPBSs with differentially expressed genes from two cell types. The method evaluates the association of the quality of footprints overlapping particular motifs and the expression of the TF.

We used limma⁴⁷ to perform between-array normalization on expression of H1-hESC, K562 and GM12878 cells and obtain fold change estimates. Then, we retrieved all non-redundant PFMs from Jaspar in which gene symbol is a perfect match with genes present in the array platform. This leads us to 143 PFMs (see **Supplementary Datasets 2b-d**). We applied a genome-wide motif matching using these PFMs.

Afterwards, we evaluated the FLR¹² score, TC¹⁵ and FS¹⁵ for the footprints of each evaluated method, which intersects with MPBSs of a particular motif. We only considered the footprints within DHSs that are in common between the cell type pair being evaluated, as described in Yardımcı et al.¹². We expect that TFs expressed in cell type A would present higher values regarding these metrics (FLR, TC and FS) with DNase-seq from cell type A in comparison with these metrics evaluated with DNase-seq from cell type B, and vice-versa. We used a two-sample Kolmogorov-Smirnov (KS) test to assess the difference between each metrics' distribution between the two cell types being evaluated. The KS statistic, which varies from 0 to 1, is used to indicate the difference between two distributions; higher values indicate higher differences. As the KS score do not indicate the direction of the changes in distribution, we obtained a signed version by multiplying KS statistic by -1 in cases where the median of A < median of B. We calculate the Spearman correlation between the signed KS test statistic and the expression fold change for each TF (see **Supplementary Fig. 1 and 2**). Positive values indicate an association between expression of TFs and quality of footprint predictions. We will call this correlation "FLR-Exp". Results for FLR-Exp analysis are summarized in **Supplementary Dataset 2a**.

Protection score. We propose a measure to detect TF-specific footprint protection for a given DNase-seq experiment and MPBSs of a given motif/TF. As previously indicated in Sung et al.⁷, fewer DNase-seq cuts (protection) surrounding the binding site characterizes TFs with shorter binding times. More formally, the protection score for a set of *MPBS* is defined as:

$$PROT_{MPBS} = \sum_{i=1}^N \frac{(N_{R,i} - n_{C,i}) + (n_{L,i} - n_{C,i})}{2N}$$

663

664

where $MPBS = \{MPBS_1, \dots, MPBS_N\}$ is set of binding sites for a given motif, $MPBS_i = [m_i, n_i]$ is the genomic location of the i -th binding site and $n_{C,i}$, $n_{L,i}$ and $n_{R,i}$ are the number of DNase-seq reads in the binding site, upstream and downstream flanking positions, respectively (see **equation (2)** for details).

666

667

668

669

670

671

672

673

674

675

676

In short, the protection score indicates the average difference of DNase-seq counts in the flanking region and the DNase-seq counts within the MPBS. Positive values will indicate protection in the flanking regions, while values close to zero or negative indicate no protection. The protection score is a similar statistic as the FS¹⁵. The main difference is that the FS score measures the ratio between reads in flanking versus binding sites, while the protection score measures the difference. Finally, since we are interested in using the protection score as a measure of quality for a given TF and set of footprint predictions, we only evaluate MPBSs overlapping with footprints for a given cell type. The DNase-seq count values are previously corrected for DHS sequence bias and coverage differences. Results for protection scores are provided in **Supplementary Dataset 1**.

677

678

679

680

681

682

Statistical methods. The non-parametric Friedman-Nemenyi hypothesis test⁴⁸ was used to compare the AUC and AUPR of the methods regarding all data set combinations (TFs versus cell types). Such test provides a rank of the methods as well as the statistical significance of whether a particular method was outperformed. All correlations are based on Spearman values. All reported p -values have been corrected with the Benjamini and Hochberg method⁴⁹.

683

684

685

686

687

Code Availability. Software, custom code, benchmarking data, DNase-seq sequence bias estimates and further graphical results are available at www.costalab.org/hint-bc. The HINT, HINT-BC and HINT-BCN softwares can be directly accessed through the regulatory genomics toolbox website at www.regulatory-genomics.org/hint/.

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

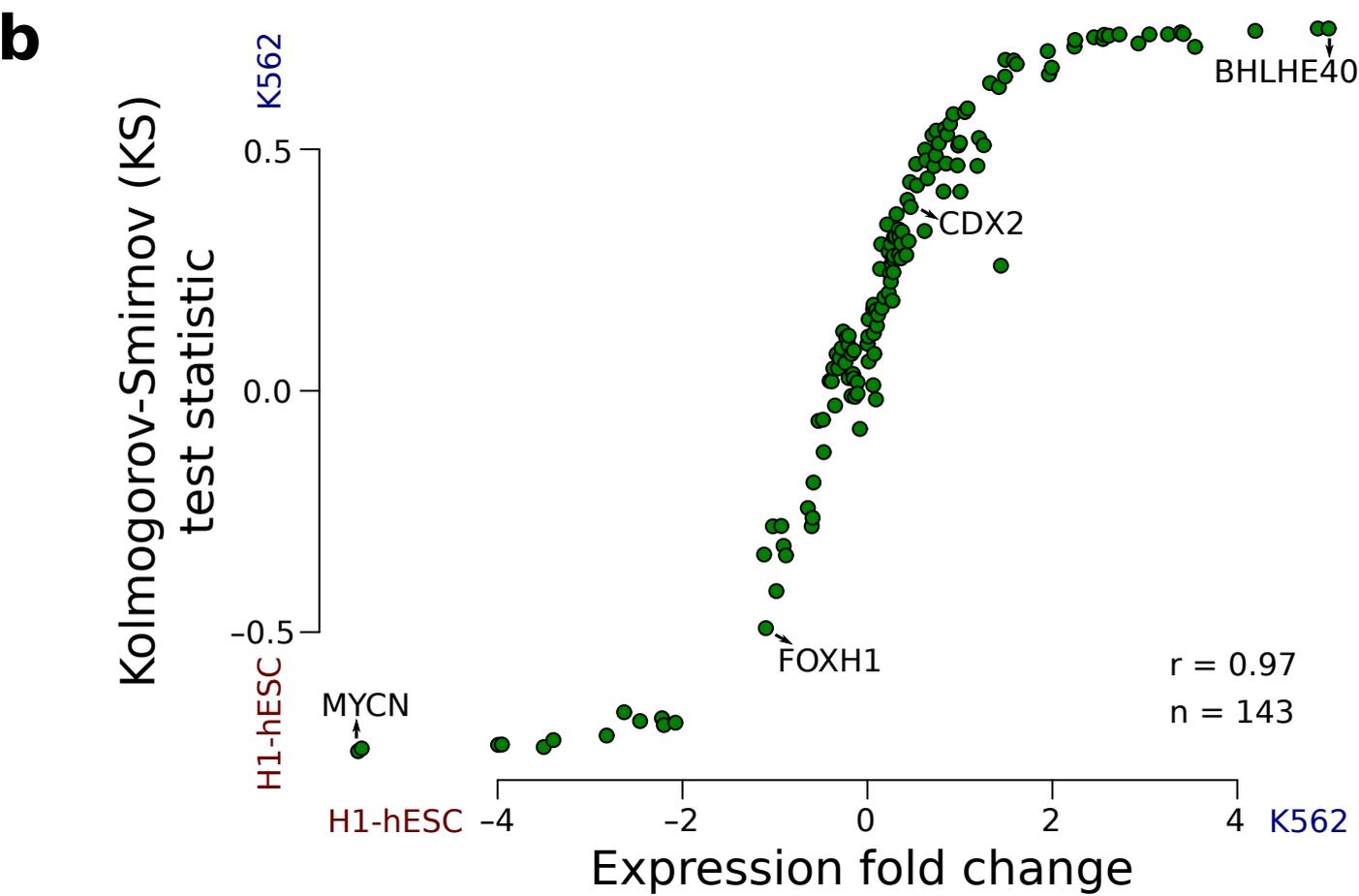
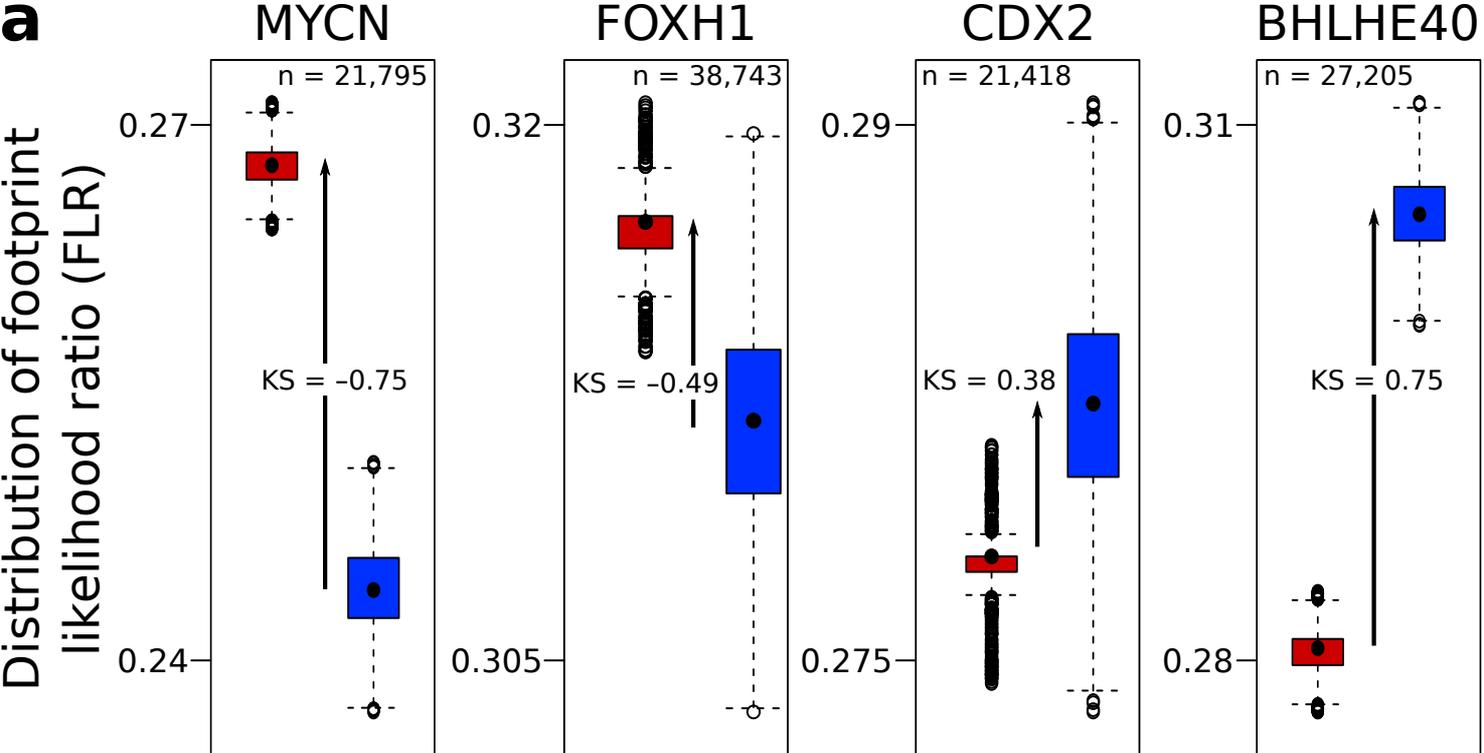
720

721

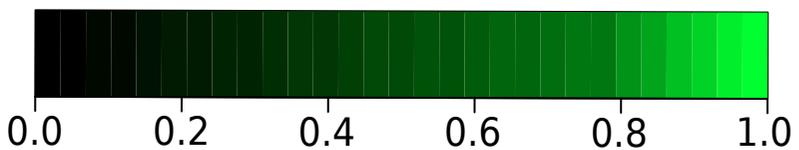
- ²⁶ Lazarovici, A. et al. Probing DNA shape and methylation state on a genomic scale with DNase I. *Proc. Natl. Acad. Sci. USA* **110**, 6376–6381 (2013).
- ²⁷ Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
- ²⁸ Zhang, Y. et al. Model-based analysis of ChIP-seq (MACS). *Genome Biol.* **9**, R137+ (2008).
- ²⁹ Yu, J. et al. An integrated network of androgen receptor, polycomb and TMPRSS2-ERG gene fusions in prostate cancer progression. *Cancer Cell* **17**, 443–454 (2010).
- ³⁰ Guertin, M.J., Zhang, X., Coonrod, S.A. & Hager, G.L. Transient estrogen receptor binding and p300 redistribution support a squelching mechanism for estradiol-repressed genes. *Mol. Endocrinol.* **28**, 1522–1533 (2014).
- ³¹ John, S. et al. Chromatin accessibility pre-determines glucocorticoid receptor binding patterns. *Nat. Genet.* **43**, 264–268 (2011).
- ³² Mathelier, A. et al. JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **42**, D142–D147 (2014).
- ³³ Robasky, K. & Bulyk, M. L. UniPROBE, update 2011: expanded content and search tools in the online database of protein-binding microarray data on protein-DNA interactions. *Nucleic Acids Res.* **39**, D124–D128 (2011).
- ³⁴ Matys, V. et al. TRANSFAC and its module TRANSCOMP: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **34**, D108–D110 (2006).
- ³⁵ Boyle, A.P., Guinney, J., Crawford, G.E. & Fury, T.S. F-seq: a feature density estimator for high-throughput sequence tags. *Bioinformatics* **24**, 2537–2538 (2008).
- ³⁶ Hesselberth, J. R. et al. Global mapping of protein-DNA interactions in vivo by digital genomic footprinting. *Nat. Methods* **6**, 283–289 (2009).
- ³⁷ Sabo, P. J. et al. Discovery of functional noncoding elements by digital analysis of chromatin structure. *Proc. Natl. Acad. Sci. USA* **101**, 16837–16842 (2004).
- ³⁸ Madden, H.H. Comments on the Savitzky-Golay convolution method for least-squares fit smoothing and differentiation of digital data. *Anal. Chem.* **50**, 1383–1386 (1978).
- ³⁹ Siepel, A. et al. Evolutionarily conserved elements in vertebrate, insect, worm and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
- ⁴⁰ Hubbard, T. et al. The ensembl genome database project. *Nucleic Acids Res.* **30**, 38–41 (2002).
- ⁴¹ Grant, C.E., Bailey, T.L. & Noble, W.S. FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27**, 1017–1018 (2011).
- ⁴² Stormo, G.D. DNA binding sites: representation and discovery. *Bioinformatics* **16**, 16–23 (2000).
- ⁴³ Cock, P.J.A. et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- ⁴⁴ Korhonen, J., Martinmäki, P., Pizzi, C., Rastas, P. & Ukkonen, E. MOODS: fast search for position weight matrix matches in DNA sequences. *Bioinformatics* **25**, 3181–3182 (2009).
- ⁴⁵ Wilczynski, B., Dojer, N., Patelak, M. & Tiuryn, J. Finding evolutionarily conserved cis-regulatory modules with a universal set of motifs. *BMC Bioinform.* **10**, 82+ (2009).
- ⁴⁶ Fawcett, T. An introduction to ROC analysis. *Pattern Recogn. Lett.* **27**, 861–874 (2006).
- ⁴⁷ Ritchie, M. E. et al. *limma* powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
- ⁴⁸ Demšar, J. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.* **7**, 1–30 (2006).

727
728
729

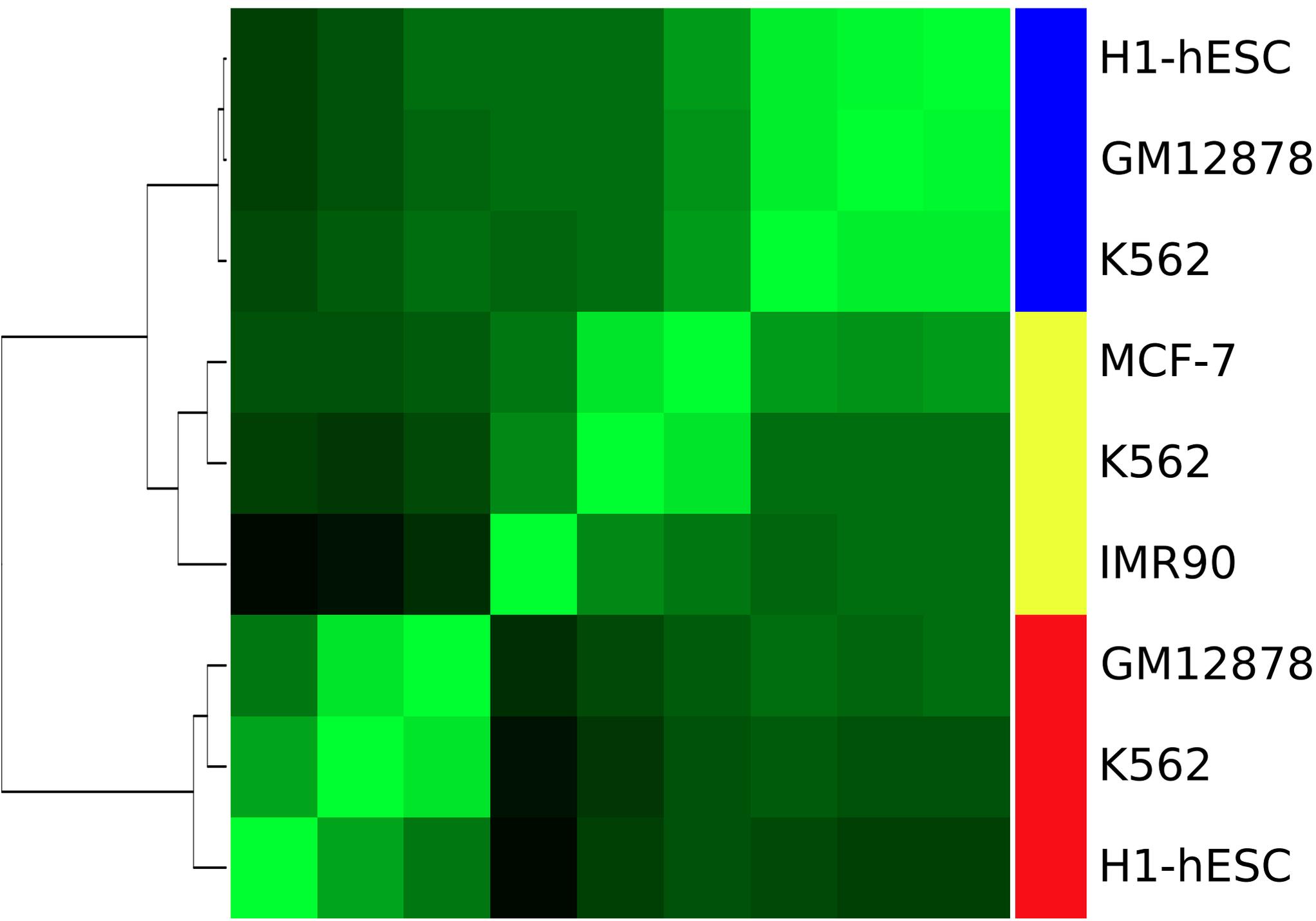
⁴⁹ Benjamini, Y. & Hochberg, Y. et al. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. R. Stat. Soc.* **57**, 289–300 (1995).

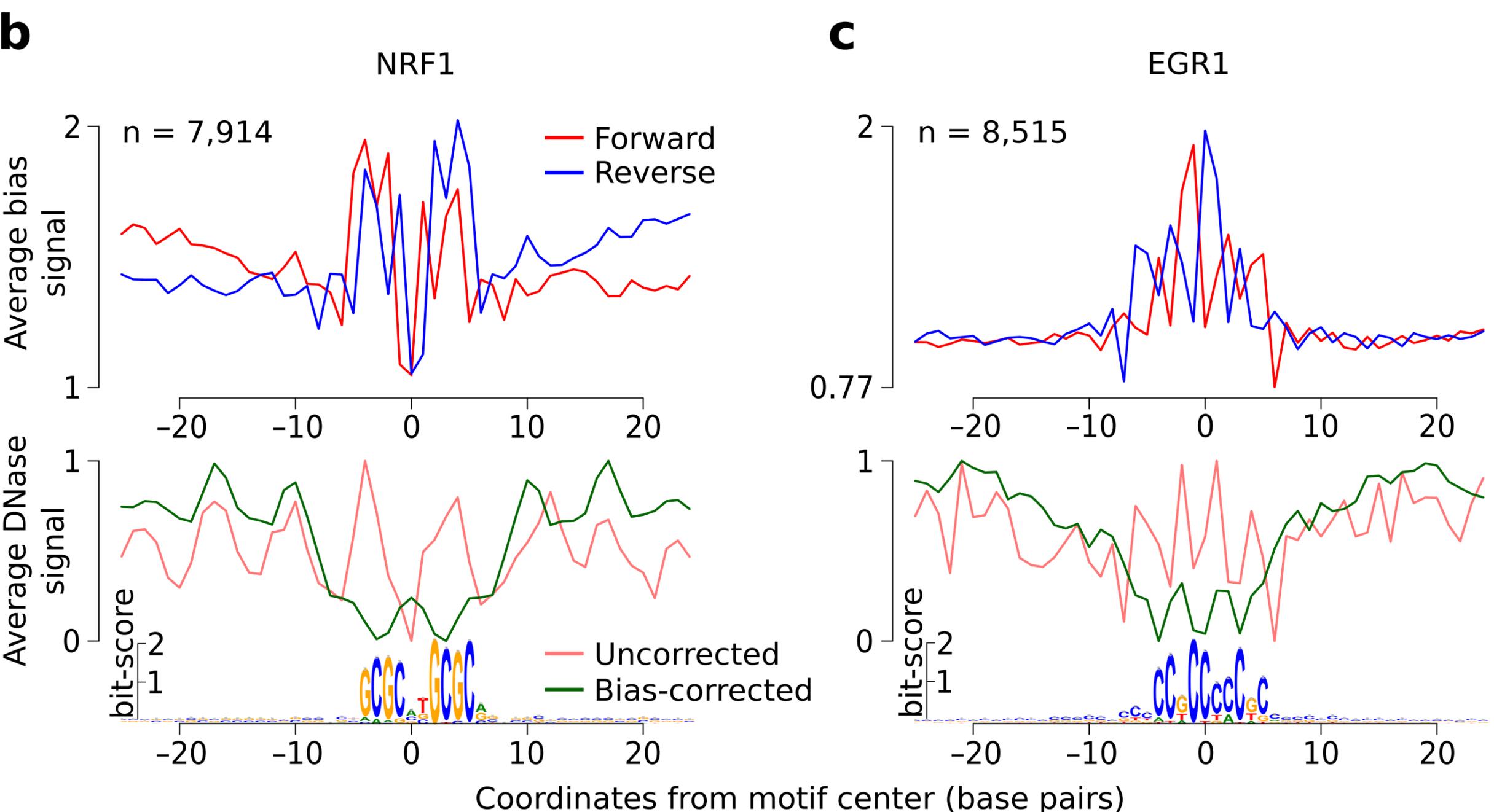
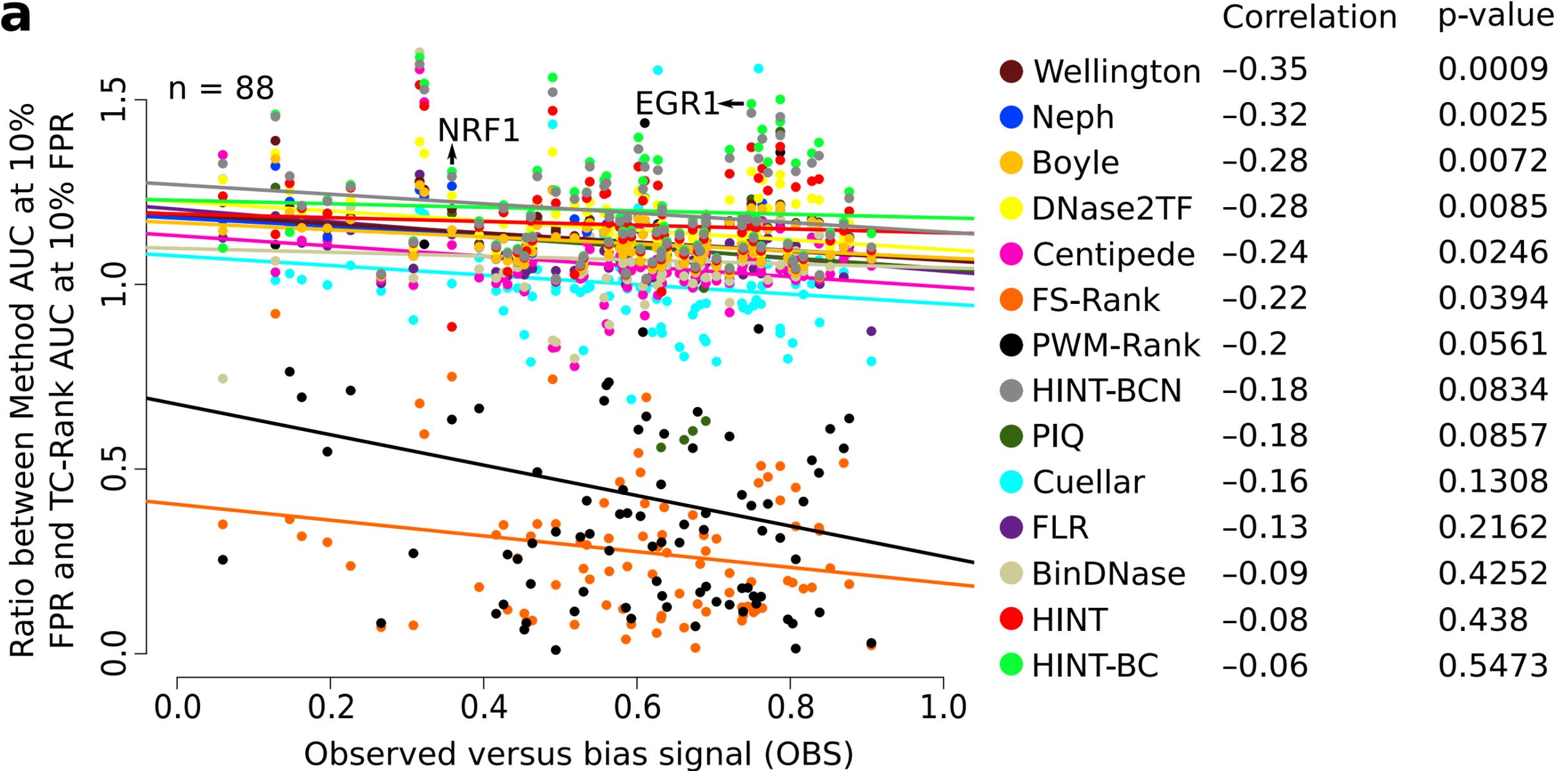


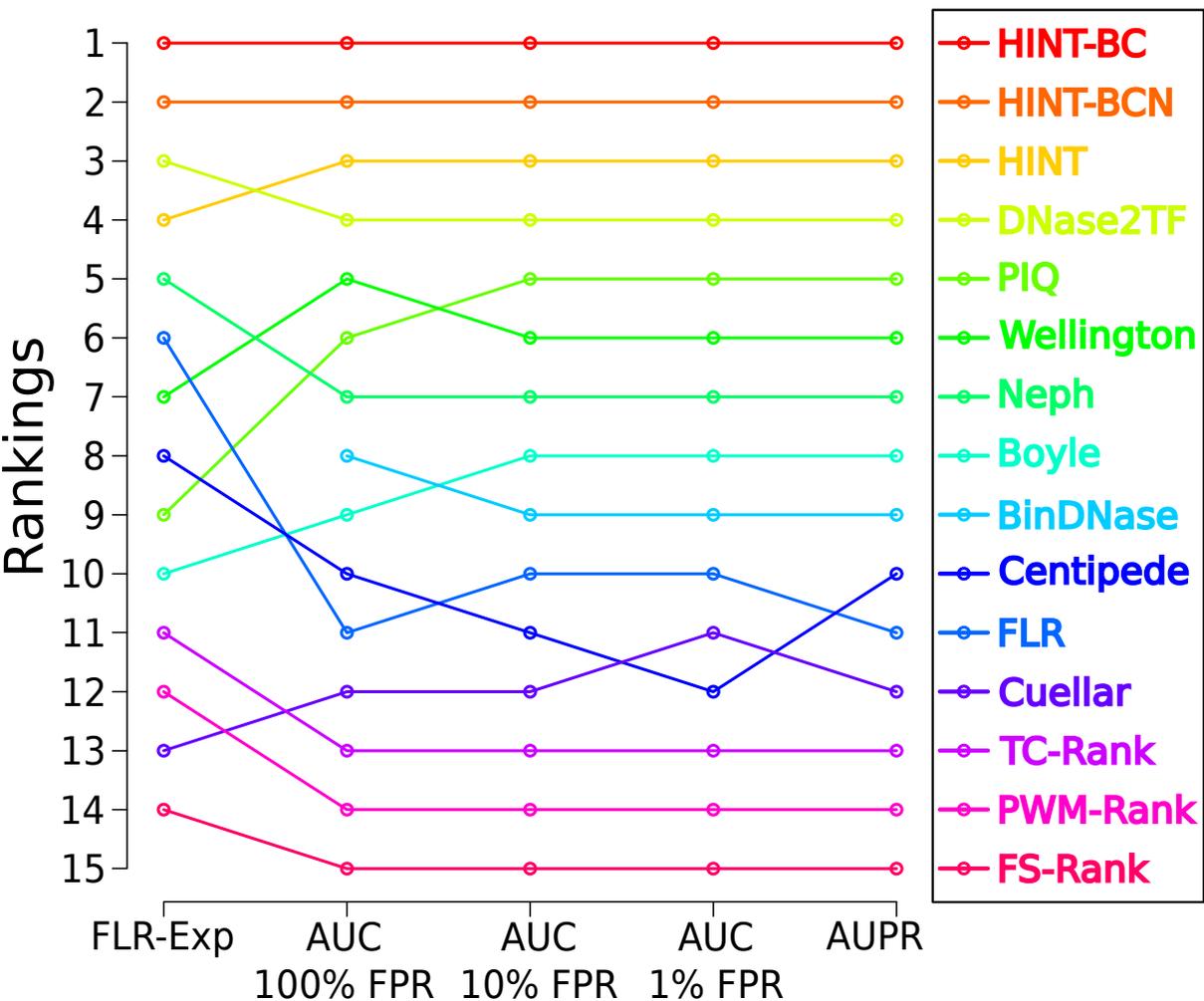
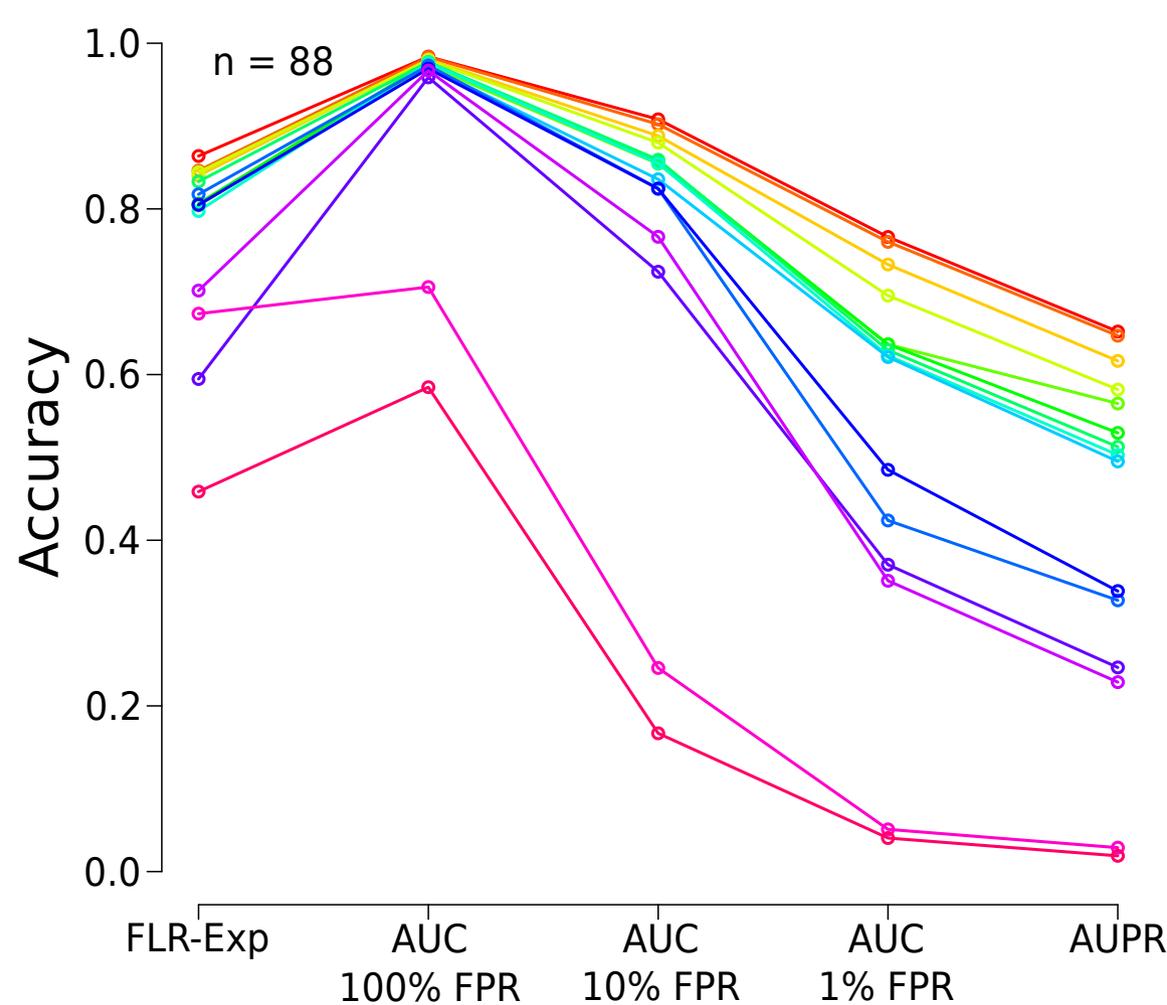
Spearman correlation

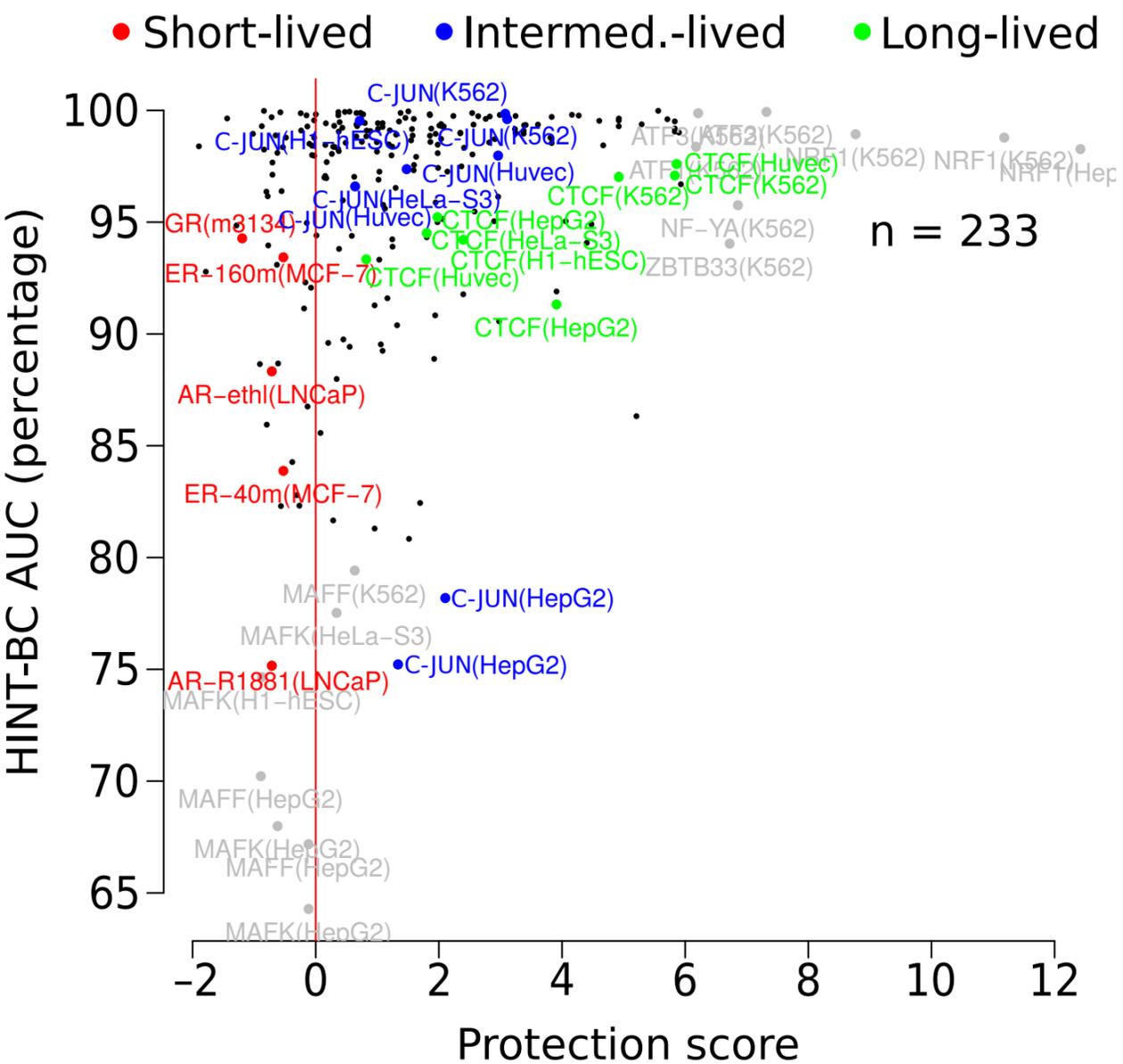
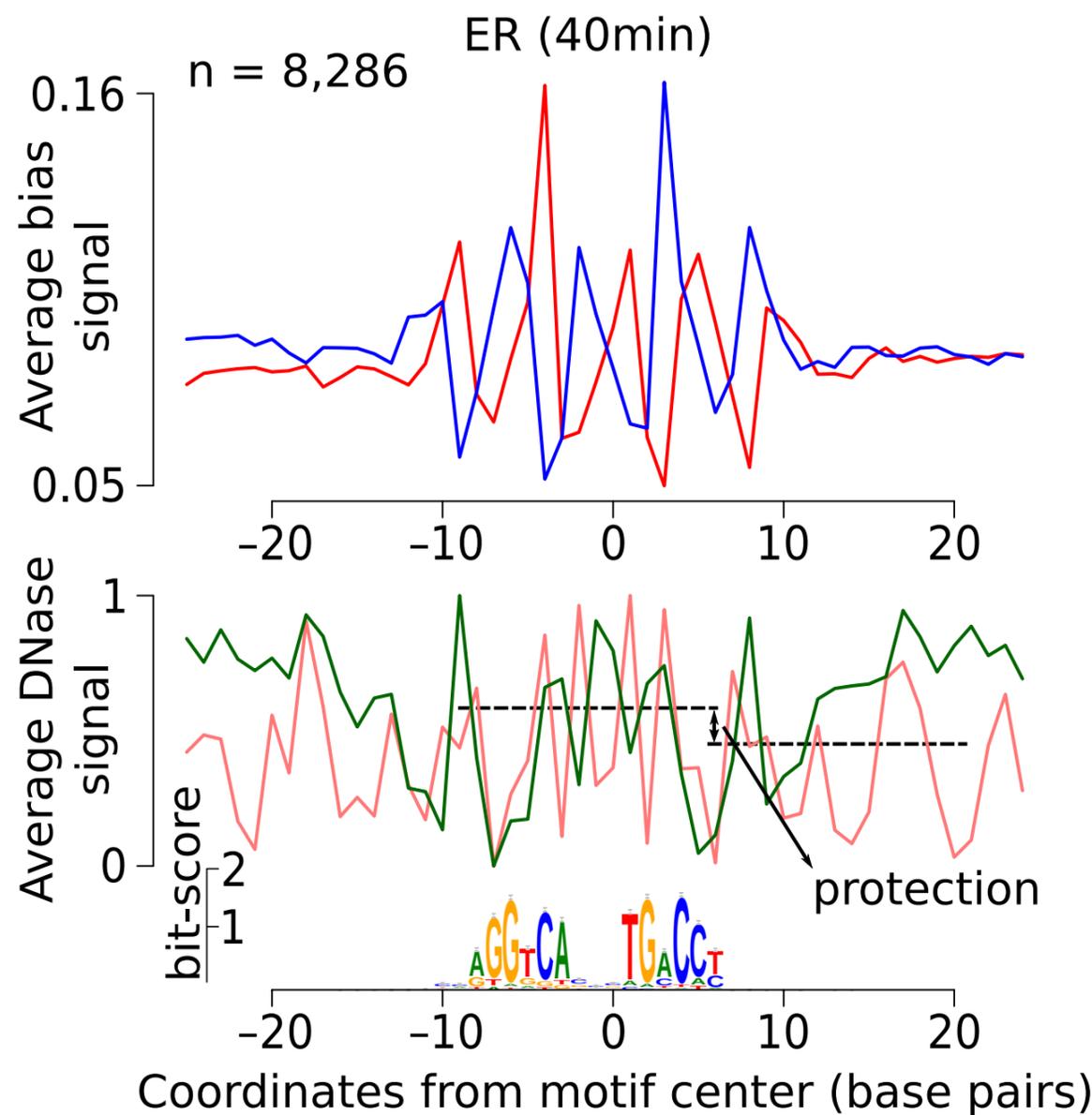
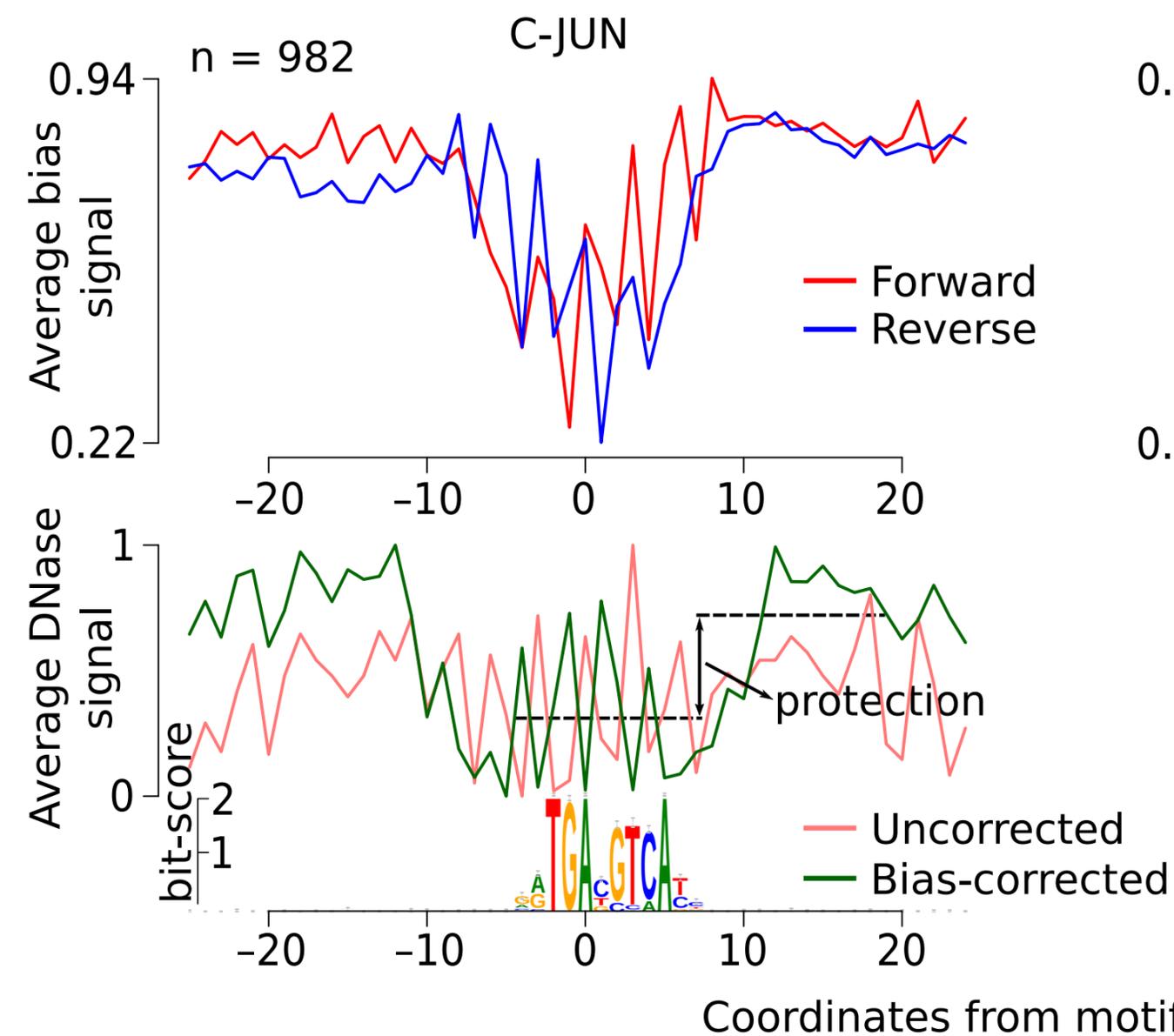


- Double-hit protocol
- Naked DNA
- Single-hit protocol





a**b**

a**b****c****d**