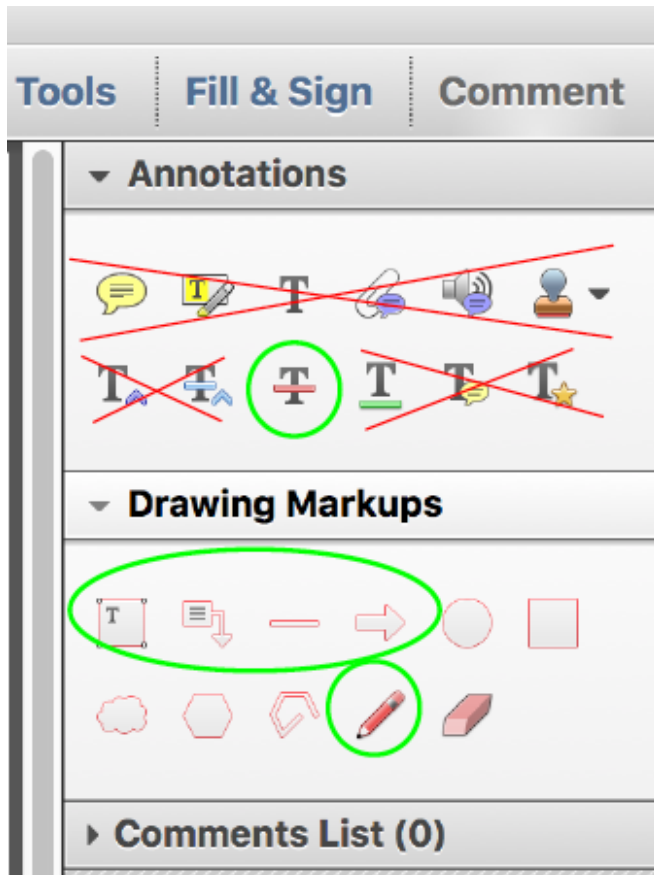


## Instructions for reviewing your article proof

You will need to use Adobe Reader version 7 or above to make comments on this proof, available free from <http://get.adobe.com/reader/>.

*Responding to the author query form:* Please address each query on the query list, on the proof or in a separate query response list.



*Editing the proof:* Edit the manuscript as necessary, using only the circled tools in the Drawing or Annotation menu.

**IMPORTANT:** All edits must be visible in full on a printed page, or they will be lost in production. Do not use any tool that does not show your changes directly on a printed page.

Use Text boxes to add text changes and the drawing tools to indicate insert points or graphics changes. Extensive directions should be addressed in a text box. While most figure edits should be marked on the page, extensive visual

changes to figures (e.g., adding scale bars or changes to data) should be accompanied by a new figure file with a written explanation of the changes.

Special characters can be inserted into text boxes by pasting from a word processing document or by copying from the list below.

α β χ δ ε φ γ η ι φ κ λ μ ν ο π θ ρ σ τ υ π ω ξ ψ ζ

A B Σ Δ Ε Φ Γ Η Ι Θ Κ Λ Μ Ν Ο Π Θ Ρ Σ Τ Υ ς Ω Ξ Ψ Ζ

Å Δ ≥ ≤ ≠ × ± 1° 3' ↑ ↓ → ←

# Analysis of computational footprinting methods for DNase sequencing experiments

■ Eduardo G Gusmao<sup>1,2</sup>, Manuel Allhoff<sup>1,3</sup>, Martin Zenke<sup>1,2</sup> & Ivan G Costa<sup>1-3</sup>✉

**DNase-seq allows nucleotide-level identification of transcription factor binding sites on the basis of a computational search of footprint-like DNase I cleavage patterns on the DNA. Frequently in high-throughput methods, experimental artifacts such as DNase I cleavage bias affect the computational analysis of DNase-seq experiments. Here we performed a comprehensive and systematic study on the performance of computational footprinting methods. We evaluated ten footprinting methods in 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 outperformed the other evaluated methods and that correcting the DNase-seq signal for experimental artifacts significantly improved the accuracy of computational footprints. We also propose a score that can be used to detect footprints arising from transcription factors with potentially short residence times.**

Next-generation sequencing combined with genome-wide mapping techniques such as DNase-seq has contributed greatly to our understanding of gene regulation and chromatin dynamics<sup>1-3</sup>. DNase-seq allows for nucleotide-level identification of transcription factor binding sites (TFBSs). This can be done via a computational search of footprint-like regions with low numbers of DNase I cuts surrounded by regions with high numbers of cuts<sup>2,3</sup>. A number of computational footprinting methods have been proposed in past years<sup>4-13</sup>. Among other applications, these methods allow the delineation of the human regulatory lexicon with millions of TFBSs over distinct cell types<sup>4</sup>, the detection of uncharacterized transcription factor (TF) motifs indicating putative regulatory elements<sup>4</sup> and the study of conservation of regulatory regions across different species<sup>14</sup>.

Next-generation sequencing-based data are significantly affected by artifacts, which are inherent to the experimental protocols used<sup>15-17</sup>. An example is the DNase I sequence cleavage bias, which is due to the different binding affinities of DNase I toward specific DNA sequences. He *et al.*<sup>15</sup> showed that sequence cleavage bias around TFBSs strongly affects the performance of a computational footprinting method<sup>4,15</sup> (footprint score (FS))

in a TF-specific manner. They also indicated several TFs, such as nuclear receptors and *de novo* motifs found via computational

**Q1. Please remove affiliation number 1 from author "Martin Zenke". This author should have only the affiliation number 2.**

comes than ranking by FS (number of DNase-seq reads inside and around a motif-predicted binding site (MPBS)). Another experimental aspect affecting the computational analysis of TF binding. Sung *et al.*<sup>7</sup> showed lower DNase I cleavage-protection 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. Although both of the aforementioned studies<sup>7,15</sup> show the challenges imposed by cleavage bias and residence time, there have been a few attempts<sup>7,12,15</sup> to address these challenges computationally.

**Please change "improved" by "improves".**

There is no well-defined gold standard for the evaluation of footprinting methods. All work so far has used chromatin immunoprecipitation followed by sequencing (ChIP-seq) of TFs in conjunction with motif-based predictions as ground truth. In short, MPBSs supported by ChIP-seq peaks are positive examples (true TFBSs), whereas MPBSs without ChIP-seq support are negative examples (false TFBSs)<sup>10</sup>. This evaluation requires TF ChIP-seq experiments to be carried out on the very same cells as the DNase-seq experiment and has a few caveats. First, TF ChIP-seq peaks are also observed in indirect binding events<sup>4,7,12,18</sup>. Second, the signal generated after TF ChIP-seq computational processing has a lower spatial resolution than DNase-seq signal. Therefore, false TFBSs might be regarded as true TFBSs if they are in close proximity to a real TFBS of the same TF<sup>15,17</sup>. Recently, Yardimci *et al.*<sup>12</sup> indicated that footprint quality scores, as measured by the footprint likelihood ratio (FLR), are significantly higher in cells where the TF in question is expressed. This observation indicates that comparing changes in the expression and quality of footprints in pairs of cells could provide an alternative footprint evaluation measure. Finally, with the exception of a few studies<sup>8,11-13</sup>, comparative analyses evaluating footprinting methods have been

<sup>1</sup>IZKF Computational Biology Research Group, RWTH Aachen University Medical School, Aachen, Germany. <sup>2</sup>Department of Cell Biology, Institute of Biomedical Engineering, RWTH Aachen University Medical School, Aachen, Germany. <sup>3</sup>Aachen Institute for Advanced Study in Computational Engineering Science (AICES), RWTH Aachen University, Aachen, Germany. Correspondence should be addressed to I.G.C. (ivan.costa@rwth-aachen.de).

Better to use: "PIQ<sup><ref 11></sup> (protein interaction quantification)" to follow the same format as HINT (in the line above). Please change <ref 11> by the superscript citation.

tational footprinting methods—Neph<sup>4</sup>, Boyle<sup>5</sup>, Wellington<sup>6</sup>, DNase2TF<sup>7</sup>, HINT<sup>8</sup> (HMM-based identification of TF footprints), CENTIPEDE<sup>9</sup>, Cuellar<sup>10</sup>, protein interaction quantification<sup>11</sup> (PIQ), FLR<sup>12</sup> and BinDNase<sup>13</sup>—with respect to their accuracy in recovering TFBSs supported by 88 ChIP-seq TF experiments in two cell types (H1-hESC and K562) with area under the receiver operating characteristic curves (AUCs) and area under the precision-recall curves (AUPRs). We also propose the FLR-Exp methodology, which associates the FLR<sup>12</sup> scores for footprints in cell-type pairs with the fold-change expression of the TFs associated with the footprints. This analysis is based on the comparison of footprints and expression of 143 TFs in H1-hESC, K562 and GM12878 cells. We also evaluated approaches for ranking footprints, strategies for dealing with DNase-seq experimental artifacts and the effect of TF residence time on footprint predictions.

## RESULTS

### Computational genomic footprinting methods

Computational footprinting methods can be broadly categorized as either segmentation<sup>4–8</sup> or site-centric<sup>9–13</sup> methods. Several segmentation methods use a window search to scan DNase-seq genomic profiles with a footprint-like shape—short regions with low DNase-seq digestion between short regions with high DNase-seq digestion (Neph<sup>4</sup>, Wellington<sup>6</sup> and DNase2TF<sup>7</sup>). Another family of segmentation methods is based on hidden Markov models, in which the hidden states model distinct levels of DNase-seq cleavage activity around footprints (Boyle<sup>5</sup> and HINT<sup>8</sup>). Site-centric methods analyze DNase-seq profiles around MPBSs and classify these sites as either bound or unbound. Most site-centric methods are based on unsupervised statistical methods such as mixture models (FLR<sup>12</sup>), Bayesian mixture models (CENTIPEDE<sup>9</sup>) and a combination of Gaussian process and expectation propagation (PIQ<sup>11</sup>). An alternative site-centric approach proposed by Cuellar *et al.*<sup>10</sup> uses DNase-seq profiles as the prior distribution for the detection of MPBSs. BinDNase is a supervised site-centric method based on logistic regression<sup>13</sup>. We also evaluated simple statistics as baseline methods, ranking MPBSs by position weight matrix (PWM-Rank) bit-score<sup>10</sup>, by the ratio of the number of DNase-seq reads inside an MPBS to that around it (FS-Rank)<sup>4,15</sup> and by the number of DNase-seq reads around an MPBS (TC-Rank)<sup>10,15</sup>.

There are several other relevant characteristics for computational footprinting methods. A few methods allow the inclusion of additional genomic and/or experimental evidence such as conservation scores<sup>9</sup>, distance to transcription start sites<sup>9</sup> and histone modifications<sup>8–10</sup>. Only PIQ<sup>11</sup> supports the analysis of several DNase-seq data sets (i.e., experiments with replicates or time series). Another important feature is the correction of DNase-seq experimental artifacts, which is supported only by DNase2TF<sup>7</sup>, HINT<sup>8</sup> variants (HINT-BC and HINT-BCN) and FLR<sup>9</sup>. Whereas HINT-BC, HINT-BCN and DNase2TF use experimental bias statistics to pre-process DNase-seq profiles, FLR builds a ‘cleavage bias’ model within the mixture model in

a TF-specific manner. Most methods use base pair DNase-seq resolution as primary input<sup>4–9,11–13</sup>. One exception is Cuellar<sup>10</sup>, which is based on smoothed DNase-seq signals from 150-bp windows. Smoothing of base pair-resolution profiles is performed by PIQ via the use of Gaussian process models<sup>11</sup>. BinDNase uses a greedy backward feature-selection approach that merges read counts of neighboring genomic positions<sup>13</sup>. Footprinting methods also provide statistics that can be used to rank footprint predictions. Wellington<sup>6</sup> and DNase2TF<sup>7</sup> use read count statistics to provide *P* values for each footprint. Several site-centric approaches provide either probabilities (BinDNase<sup>13</sup>, CENTIPEDE<sup>9</sup> and PIQ<sup>11</sup>) or log-odds scores (FLR<sup>12</sup>) of footprints. Other methods use statistics such as FS (Neph<sup>4</sup>), position weight matrix scores (Cuellar<sup>10</sup>) and TC (HINT<sup>8</sup>) to rank predicted footprints.

The availability, usability and scalability of the software tools implementing the methods are also important features. Neph<sup>4</sup>, HINT<sup>8</sup>, PIQ<sup>11</sup> and Wellington<sup>6</sup> provide tutorials and software for running experiments with few command-line calls. Of those, only HINT<sup>8</sup>, PIQ<sup>11</sup> and Wellington<sup>6</sup> natively support standard genomic formats as input. The site-centric methods Cuellar<sup>10</sup>, BinDNase<sup>13</sup>, CENTIPEDE<sup>9</sup> and FLR<sup>12</sup> require a single execution and input data source per TF and cell, whereas segmentation methods require an execution per cell only. These site-centric methods had computational demands 5 times (FLR and Cuellar) to 50 times (BinDNase and CENTIPEDE) higher than those of the slowest segmentation method (Wellington) in our analysis (Supplementary Table 1). The main method features are summarized in Table 1 and described in the Online Methods.

### Association of TF expression with footprint quality

Yardimci *et al.*<sup>12</sup> indicated that the FLRs of candidate footprints are significantly higher in cells where the relevant TF is being expressed. We expanded on this idea by evaluating whether differences in the FLR score distribution of footprints overlapping with MPBSs on a pair of cell types are proportional to differences in the expression of the respective TFs (Fig. 1a). We observed high average correlation values for the majority of evaluated methods ( $r = 0.79$ ) and extremely high correlation values ( $r > 0.9$ ) for the top-performing methods in comparisons between pairs of H1-hESC, K562 and GM12878 cells (Fig. 1b and Supplementary Fig. 1). We also evaluated the use of the TC and FS metrics as quality scores instead of the FLR and found that they had lower average correlation values (TC,  $r = 0.35$ ; FS,  $r = 0.73$ ; Supplementary Fig. 2). We therefore opted to use the FLR as the quality measure for footprints for this evaluation procedure. We used the correlation between FLR score difference and expression fold change, which we refer to as FLR-Exp, to rank footprinting methods; higher FLR-Exp values indicate better performance. The FLR-Exp evaluation methodology requires only expression data and is therefore more generally applicable than TF ChIP-seq-based evaluation. However, unlike TF ChIP-seq evaluation, the FLR-Exp approach cannot evaluate footprint predictions for individual TFs.

### Impact of experimental artifacts

To understand the nature of artifacts in DNase-seq experiments, we analyzed the sequence bias estimates of all 61 ENCODE tier 1 and 2 DNase-seq data sets (Supplementary Table 2). These

### Table 1 | Overview of methods

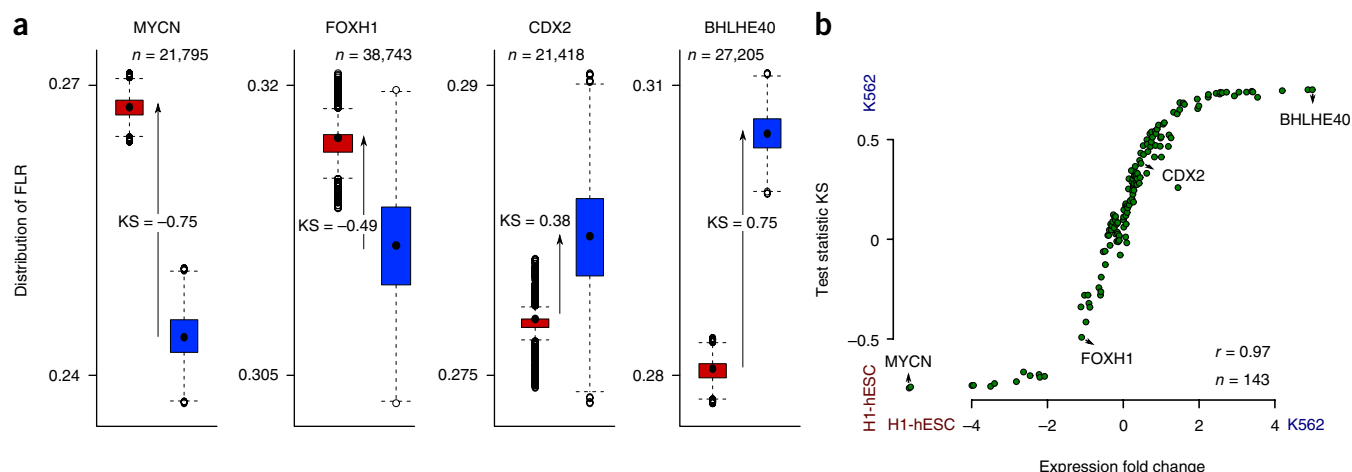
Name	Type	Algorithm	Bias correction	Resolution	Score	+	-	Notes
BinDNase	SC	Logistic regression	None	Base pair				
Boyle	SEG	HMM	None	Base pair				
CENTIPEDE	SC	Bayesian mixture model	None	Base pair				
Cuellar	SC	Weighted motif match	None	Sliding window	PWM score	+	-	
DNase2TF	SEG	Sliding window	4-mer (DHS sequence bias)	Base pair	P 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				Integrates histones
Neph	SEG	Sliding window	None	Base pair				
PIQ	SEG	GP, expectation propagation	None	Base pair				Supports replicates, time series
Wellington	SEG	Sliding window	None	Base pair	P value	+	+	

Main characteristics of the evaluated methods. Under “Availability,” “+” denotes that the method is publicly available. The Boyle method is not public, but the authors provide footprint predictions for a few cells. Under “Usability,” “+” denotes that the method natively supports standard genomic files and can be executed with few commands ( $\leq 3$ ). SC, site centric; SEG, segmentation approach; HMM, hidden Markov model; PWM, position weight matrix; GP, Gaussian process.

experiments included two main DNase-seq protocols that differ in the number of DNase I digestion events necessary to generate DNA fragments (single-hit<sup>2</sup> and double-hit<sup>3</sup>). Sequence bias estimates can be defined as the ratio of observed to expected DNase-seq reads starting at the middle of a particular DNA sequence of length  $k$  ( $k$ -mer)<sup>15</sup>. We used two approaches. The DNase hypersensitive site (DHS) sequence bias approach considers the sequence bias estimates within DHSs in each DNase-seq experiment. This approach captures the DNase I cleavage, read fragmentation and sequence-complexity bias of DHSs in each DNase-seq experiment<sup>7,15</sup>. The naked DNA sequence bias approach considers the sequence bias estimates in naked DNA DNase-seq experiments<sup>12</sup>. In this case, all DNA regions are open, and therefore the sequence bias estimates capture mainly the DNase I cleavage bias<sup>12</sup> (Online Methods). Our clustering

If consistent with journal format, please change the dot here for a colon punctuation.

Next we extended the analysis by He *et al.*<sup>15</sup> to evaluate the influence of sequence bias on all evaluated footprinting methods on the basis of the AUC at a 10% false positive rate (FPR). In this analysis, the Spearman correlation is calculated for the distribution of TF accuracy for each method in comparison to the distribution of TFs' observed (uncorrected) DNase-seq signal versus the bias signal. We evaluated the correlation between

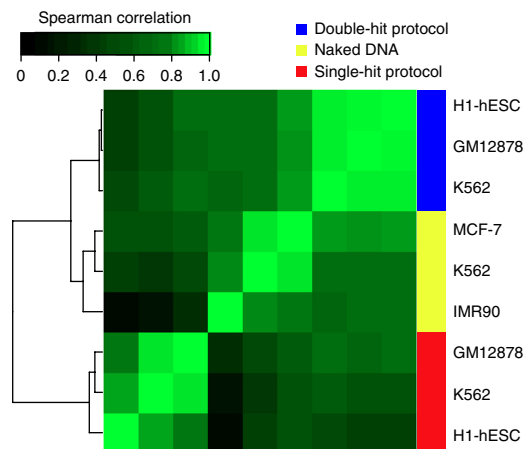


**Figure 1** | FLR-Exp evaluation metric. **(a)** FLR score distribution of footprints predicted with HINT-BC overlapping with MPBSs of selected TFs. These TFs have increased expression in K562 cells (red) compared with H1-hESC cells (blue). The Kolmogorov-Smirnov (KS) statistic quantifies the separation of the two distributions. Box plots depict distribution median values (center dots) and first and third quartiles (box edges). The whiskers represent 1.5 times the interquartile range, and external dots represent outliers (data greater than or less than 1.5 times the interquartile range). **(b)** KS statistics and expression fold change for 143 TFs in H1-hESC and K562 cells. There is a clear association between TF expression and KS statistic ( $r = 0.97$ , adjusted  $P < 10^{-10}$ ).



**Figure 2** | Clustering of bias estimates. Ward's minimum variance clustering based on the pairwise Spearman correlation coefficient ( $r$ ) from bias estimates of selected ENCODE tier 1 and naked DNA DNase-seq data. DNase-seq experiments were based on single-hit or double-hit protocols or on naked DNA.

the uncorrected and bias signals for each TF by measuring the uncorrected DNase-seq signal and the bias signal for every MPBS that overlapped a footprint from the evaluated method. Then we evaluated the Spearman correlation between the average uncorrected and bias signals. We evaluated HINT using DNase-seq signals corrected with either DHS sequence bias (HINT bias-corrected (HINT-BC)) or naked DNA sequence bias (HINT bias-corrected on naked DNase-seq (HINT-BCN)). Our analysis showed that only six out of ten evaluated methods (Wellington, Neph, Boyle, DNase2TF, CENTIPEDE and FS-Rank) presented a significant negative Spearman correlation ( $r = -0.35, -0.32, -0.28, -0.28, -0.24$  and  $-0.22$ , respectively) between their accuracy and the amount of sequence bias (Fig. 3a; adjusted  $P$  value  $< 0.05$ ). Equivalent results were also observed for the same TFs and cellular conditions analyzed by He *et al.*<sup>15</sup> (Supplementary Fig. 4). Methods explicitly using 6-mer sequence bias statistics (HINT-BC, HINT-BCN and FLR) or including smoothing (Cuellar, BinDNase and PIQ) were not significantly influenced by sequence bias. Moreover, the performance of HINT-BC was the least affected by sequence bias ( $r = -0.06$ ). Pairwise comparison of AUCs at 10% FPR for all three HINT variants (HINT-BC, HINT-BCN and HINT) indicated significant gain in all predictions with



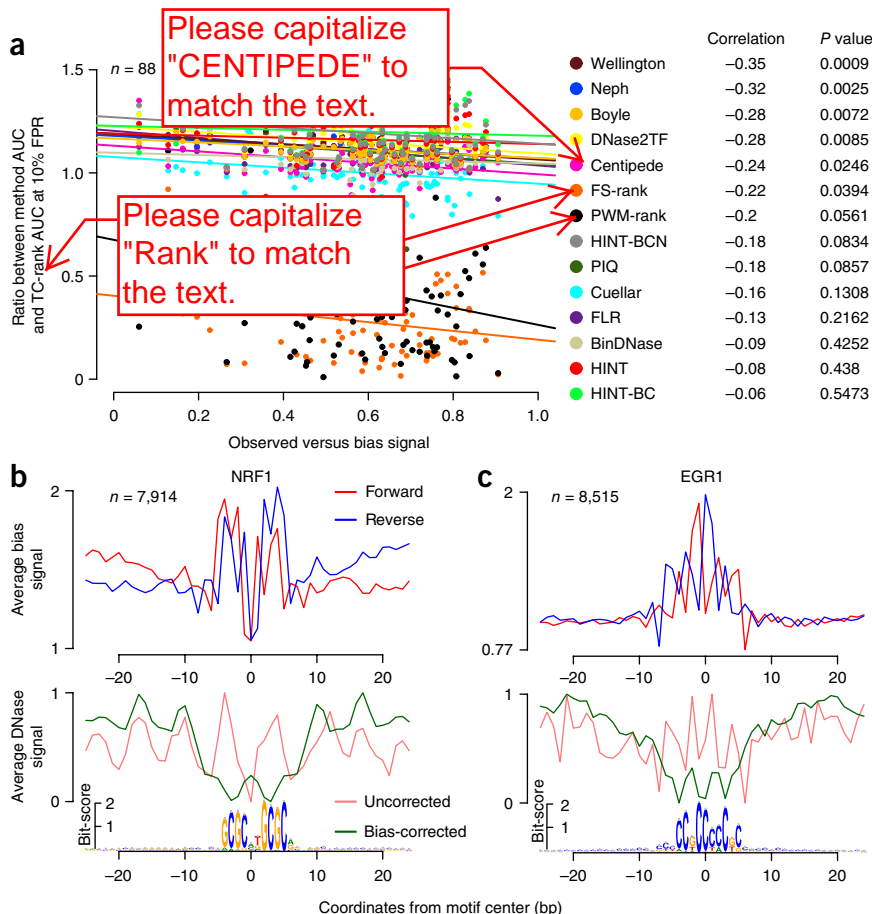
sequence bias correction (adjusted  $P < 10^{-30}$ ; Supplementary Fig. 5a). There was no significant difference between HINT-BC and HINT-BCN, but we observed a higher AUC with HINT-BC for all but seven TFs. This indicates an advantage of DHS sequence bias correction for the footprint prediction problem.

As an example, we show sequence bias estimates and corrected and uncorrected DNase-seq average profiles around TFBSs with the highest AUC gain between HINT-BC and HINT (Fig. 3b,c and Supplementary Fig. 6). The NRF1 and EGR1 DNase-seq profiles indicate that the bias-corrected signal fit their sequence affinity better than the uncorrected signal. We observed that  $k$ -mers with high DHS sequence bias had a high CG content ( $r > 0.8$  in

11 out of 12 cell types; Supplementary Fig. 7). However, there was no significant correlation between the CG content of MPBSs and either AUC values or differences in AUC among HINT-BC, HINT-BCN and HINT ( $P > 0.05$ ; Supplementary Fig. 5b).

### Comparative analysis of footprinting methods

Given its good performance<sup>10,15</sup>, we evaluated the use of the TC as the ranking strategy instead of each method's own ranking for BinDNase, CENTIPEDE, Cuellar,

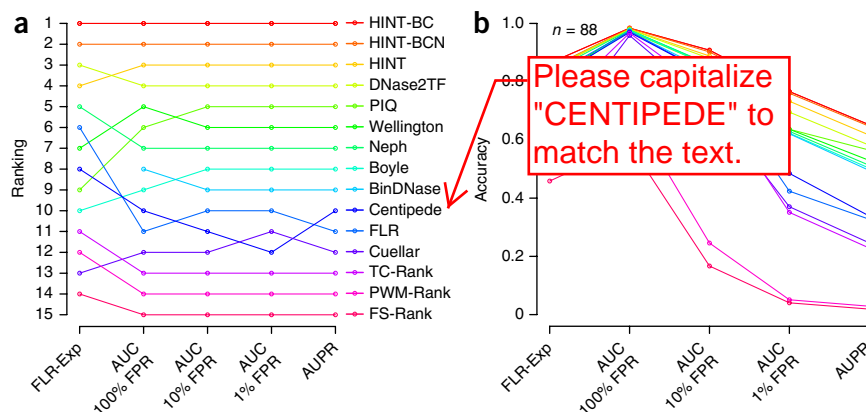


**Figure 3** | Effects of sequence bias on methods.

(a) Association between the performance of footprinting methods (relative to TC-Rank performance) and their sequence bias estimated for 88 TFs binding on H1-hESC and K562 cells. The x-axis represents the correlation between the uncorrected and bias signals (OBS). Higher OBS values indicate higher bias. The y-axis represents the ratio between the AUC at 10% FPR for each evaluated method and for the TC-Rank method; higher values indicate higher accuracy. The diagonal lines in the plot represent the linear regression's best fit.

(b,c) Average bias signal (top) and uncorrected and bias-corrected DNase-seq signals (bottom) for the TFs (b) NRF1 and (c) EGR1. DNase signals were standardized to be in the interval [0,1]. The motif logos represent all underlying DNA sequences centered on the TFBSs.

**Figure 4** | Evaluation of computational footprinting methods. (a) Average rankings for the evaluated computational footprinting methods. The rankings are given for all evaluation criteria: FLR-Exp, TF ChIP-seq-based AUC (at 100%, 10% and 1% FPR) and AUPR. (b) FLR-Exp values (as a combination of all pairwise comparisons among H1-hESC, K562 and GM12878 cells), median TF ChIP-seq-based AUC values (at 100%, 10% and 1% FPR) and median AUPR values for all evaluated methods. HINT-BC, HINT-BCN, HINT and DNase2TF were ranked as the top four methods by all evaluation metrics. All baseline methods (FS-Rank, PWM-Rank and TC-Rank) were in the bottom four positions of the ranks. Note that BinDNase could not be evaluated with the FLR-Exp because it requires ChIP-seq data for training.



DNase2TF, FLR, PIQ and Wellington. Prior to ranking by TC, site-centric methods required the definition of a minimum probability score to define active footprints. In all cases, using TC yielded higher AUC values (10% FPR) than did using the intrinsic ranking metric (**Supplementary Fig. 8**). A probability cutoff of 0.9 yielded the highest AUCs for site-centric methods, with the exception of BinDNase (highest AUC at 0.8). These parameters were used in our subsequent evaluation analyses.

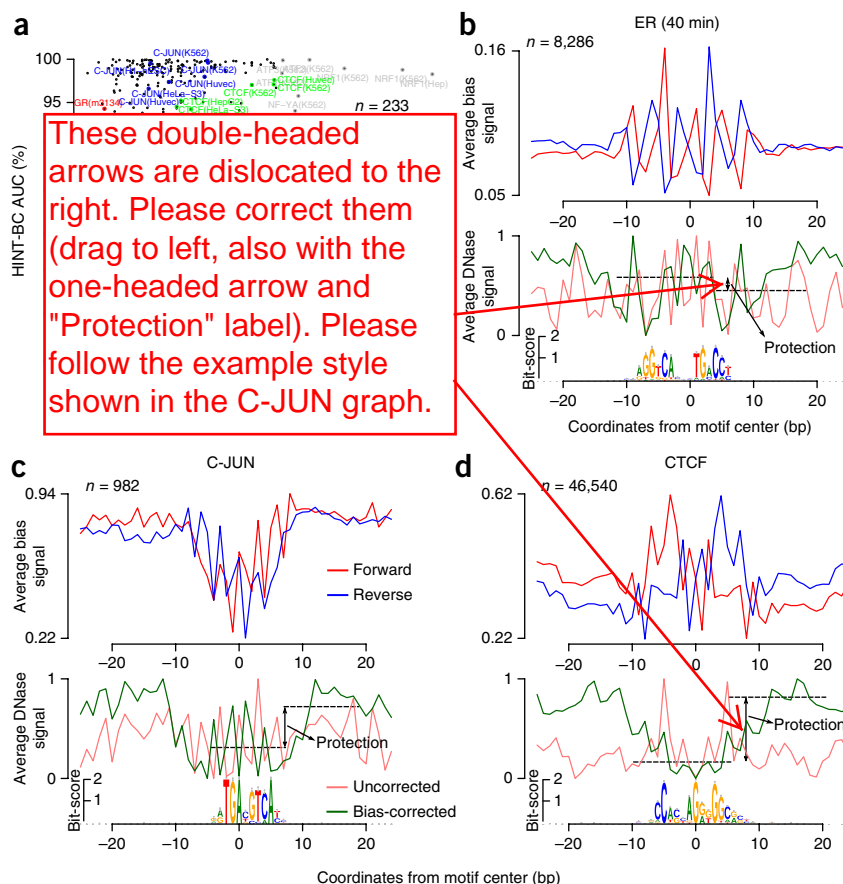
We next evaluated all the competing methods by measuring the AUC at 1%, 10% and 100% FPR using TF ChIP-seq data. AUCs at lower FPRs favored methods with higher sensitivity at the expense of specificity. We also estimated the AUPR, which

is indicated for cases with an imbalance between positive and negative examples<sup>20</sup>, and the FLR-Exp metric. All TF ChIP-seq-based metrics indicated a very similar ranking ( $r > 0.98$ ; **Fig. 4a**). There was also strong agreement between FLR-Exp and other metrics ( $r > 0.88$ ). HINT-BC had the highest FLR-Exp, AUC and AUPR values and significantly outperformed all other methods except HINT-BCN (adjusted  $P < 0.01$ ; **Supplementary Fig. 9** and **Supplementary Tables 3–6**). When we ignored HINT variants, the next top-performing method was DNase2TF, which significantly outperformed all other methods but PIQ (adjusted  $P < 0.01$ ). PIQ outperformed all of its lower-ranked competitors except for Wellington in terms of AUC (1% FPR) and AUPR (adjusted

$P < 0.01$ ). Concerning the performance of TC-Rank, we observed that its AUC values for 10% and 100% FPR were very close to those for other footprinting methods (**Fig. 4b** and **Supplementary Fig. 9**). This was not the case for values of AUC at 1% FPR or of AUPR, for which all methods but CENTIPEDE and Cuellar had significantly superior performance relative to that of TC ( $P < 0.01$ ; **Supplementary Tables 3–6**).

#### TF residence time

Despite the high average prediction values of the top-performing footprint methods, they consistently performed worst with a similar set of TFs; HINT-BC, DNase2TF



**Q6.** These labels are allowed to overlap since the position of the dots is the important information. Furthermore the TFs are being clearly explicated in the figure caption.

Average bias signal (top) and uncorrected and bias-corrected DNase-seq signals (bottom) for the TFs (b) ER (c) C-JUN and (d) CTCF. DNase-seq signals were standardized to be in the interval [0,1]. Motif logos represent all underlying DNA sequences centered on the TFBSs.

and PIQ had 89% of TFs in common in the lower quartile of the AUC at 10% FPR (**Supplementary Data Set 1**). This list included nuclear receptors, which have low residence binding times<sup>7</sup> and display a lower DNase I cleavage-protection pattern (**Supplementary Fig. 10**). To further investigate this, we proposed a statistic inspired by the concepts presented by Sung *et al.*<sup>7</sup> to detect TFs with potential short residence times, the protection score, which measures the difference between the amounts of DNase I digestion in the flanking regions and in the TFBS on bias-corrected DNase-seq signals. We used this statistic to analyze the predictive performance of different methods with TFs with distinct residence times. For this we used a comprehensive data set with 233 combinations of DNase-seq experiments and TFs (Online Methods).

We observed that TFs with known short residence times on DNA, such as nuclear receptors AR<sup>21</sup>, ER<sup>22</sup> and GR<sup>23</sup>, presented a negative protection score (**Fig. 5a**). TFs with intermediate and long residence times on DNA (C-JUN<sup>24</sup> and CTCF<sup>25</sup>, respectively) presented a positive protection score. The amount of protection was clearly reflected in the bias-corrected DNase-seq profiles (**Fig. 5b–d**). We also noted association of the protection score and the AUC in HINT-BC (**Fig. 5a**). Overall, the protection score positively correlated with the AUC values of evaluated methods such as TC ( $r = 0.19$ ) and HINT-BC ( $r = 0.26$ ) and negatively correlated ( $r = -0.49$ ) with the sequence bias (adjusted  $P < 0.05$ ). These results reinforce the concept that TFs with potentially short residence times are poorly detected via DNase-seq footprints.

## DISCUSSION

Our comparative analysis indicates the superior performance of (in decreasing order) HINT, DNase2TF and PIQ in the prediction of active TFBSs in all evaluated scenarios. Moreover, tools implementing these methods were user-friendly and had lower computational demands than other evaluated methods. Clearly, the choice of computational footprinting approach should also be based on experimental design aspects. For example, PIQ is the only method supporting analysis of replicates and time series. For studies requiring footprint predictions for *de novo* motif analysis, one should use a segmentation approach such as HINT or DNase2TF. In contrast to positive evaluations of TC-Rank in previous works<sup>10,15</sup>, we show that it has poor sensitivity as indicated by the AUC at low FPR levels. The TC statistic provided the best strategy for ranking footprint predictions from other methods.

The refined DNase-seq protocol and experimental artifacts presented by He *et al.*<sup>15</sup> and the TF binding time presented by Sung *et al.*<sup>7</sup> underscore the idea that robust *in silico* techniques are required to correct for experimental artifacts and to derive valid biological predictions. The correction of DNase-seq signal with DHS sequence bias estimates virtually removes the effects of sequence bias artifacts on computational footprinting. We have demonstrated that such correction can be performed before the execution of the computational footprinting method. It should still be noted, however, that ignoring experimental artifacts might lead to false predictions, as observed previously for predicted *de novo* motifs (**Supplementary Fig. 11**). Moreover, the simple protection score can indicate footprints of TFs with potentially short binding times. Thus, footprint predictions for TFs with low protection scores should be interpreted with caution.

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

## METHODS

Methods and any associated references are available in the [online version of the paper](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

This work was supported by the IZKF Aachen, RWTH Aachen University Medical School, Aachen, Germany (to E.G.G., M.A. and I.G.C.) and the Excellence Initiative of the German Federal and State Governments, and the German Research Foundation (grant GSC 111 to M.A. and I.G.C.).

## AUTHOR CONTRIBUTIONS

E.G.G., M.Z. and I.G.C. designed the research. E.G.G. wrote HINT program code. E.G.G., M.A. and I.G.C. analyzed data. E.G.G., M.Z. and I.G.C. wrote the manuscript.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
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).
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).
4. Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489**, 83–90 (2012).
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).
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).
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).
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).
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).
10. Cuellar-Partida, G. *et al.* Epigenetic priors for identifying active transcription factor binding sites. *Bioinformatics* **28**, 56–62 (2012).
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).

**Q2. Acknowledgements are correct.**

**Q3. Competing financial interests are correct.**

12. Yardimci, G.G., Frank, C.L., Crawford, G.E. & Ohler, U. Explicit DNase sequence bias modeling enables high-resolution transcription factor footprint detection. *Nucleic Acids Res.* **42**, 11865–11878 (2014).
13. Kähärä, J. & Lähdesmäki, H. BinDNase: a discriminatory approach for transcription factor binding prediction using DNase I hypersensitivity data. *Bioinformatics* **31**, 2852–2859 (2015).
14. Stergachis, A.B. *et al.* Conservation of trans-acting circuitry during mammalian regulatory evolution. *Nature* **515**, 365–370 (2014).
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).
16. Meyer, C.A. & Liu, X.S. Identifying and mitigating bias in next-generation sequencing methods for chromatin biology. *Nat. Rev. Genet.* **15**, 709–721 (2014).
17. Park, P.J. ChIP-seq: advantages and challenges of a maturing technology. *Nat. Rev. Genet.* **10**, 669–680 (2009).
18. Teytelman, L., Thurtle, D.M., Rine, J. & van Oudenaarden, A. Highly expressed loci are vulnerable to misleading ChIP localization of multiple unrelated proteins. *Proc. Natl. Acad. Sci. USA* **110**, 18602–18607 (2013).
19. The difficulty of a fair comparison. *Nat. Methods* **12**, 273 (2015).
20. Davis, J. & Goadrich, M. The relationship between precision-recall and ROC curves. *Proc. 23rd International Conference on Machine Learning—ICML 2006* 233–240 (2006).
21. Tewari, A.K. *et al.* Chromatin accessibility reveals insights into androgen receptor activation and transcriptional specificity. *Genome Biol.* **13**, R88 (2012).
22. Sharp, Z.D. *et al.* Estrogen-receptor- $\alpha$  exchange and chromatin dynamics are ligand- and domain-dependent. *J. Cell Sci.* **119**, 4101–4116 (2006).
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).
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).
25. Nakahashi, H. *et al.* A genome-wide map of CTCF multivalency redefines the CTCF code. *Cell Rep.* **3**, 1678–1689 (2013).

## EDITORIAL SUMMARY

**AOP:** This comparison of ten computational methods for detecting transcription factor binding sites in DNase hypersensitive regions in the genome determines which methods work consistently well, how DNase-seq experimental artifacts should be corrected for and which score is best for ranking methods.



## ONLINE METHODS

**Data.** DNase-seq aligned reads were obtained from ENCODE<sup>1</sup>. To perform the computational footprint experiments, we obtained data for H1-hESC, HeLa-S3, HepG2, HUVEC, K562, LNCaP and MCF-7 cells from Crawford's lab at Duke University (labeled "DU" in the following) and for H7-hESC, HepG2, HUVEC, K562 and m3134 cells from Stamatoyannopoulos's lab at the University of Washington (labeled "UW" in the following). We also used naked DNA (deproteinized) DNase-seq experiments with MCF-7 and K562 cells (DU)<sup>12</sup> and with IMR90 cells (UW)<sup>26</sup>. Experiments on DU cells followed the single-hit protocol, and experiments on UW cells followed the double-hit protocol. In addition, to perform the DNase-seq bias estimation clustering, we used all tier 1 and tier 2 cell types from ENCODE<sup>1</sup>. **Supplementary Table 2** presents a full description of DNase-seq data.

TF ChIP-seq enriched regions (peaks and summits) were obtained from the ENCODE analysis working group<sup>1</sup> track with the exception of the following experiments, in which the enriched regions were obtained using Bowtie 2 (ref. 27) and MACS<sup>28</sup>. AR (R1881 treatment) obtained from the GSM118811 estradiol treatment were obtained from (ref. 30). GR (dexamethasone treatment) ChIP-seq raw sequences for m3134 cells were obtained from the Sequence Read Archive (SRA) under study number SRP004871 (ref. 31). All organism-specific data (DNase-seq and ChIP-seq) are based on human genome build 37 (hg19), except the DNase-seq data for m3134 and ChIP-seq data for GR, which were based on mouse genome build 37 (mm9). Chromosome Y was removed from all analyses. Expression profiles of H1-hESC, K562 and GM12878 cells were obtained from ENCODE<sup>1</sup> (GSE12760 and GSE14863).

Please change all occurrences of this name to "Yardimci" to match the main text and because the dotless i does not work here.

Please make all indicator functions bold. The indicator function is a bold "1". If not possible, please use an alternative notation for the indicator function such as "Ind".

Please replace this N by the symbol for the "set of natural numbers" (The hollow N). I am not able to paste it here. If not possible, please replace this by: "... and each integer xi is the number of ..." (maintaining the "i" as a subscript of "x").

ensemble/encode/supplementary/integration\_data\_jan2011/byDataType/footprints/jan2011/de.novo.pwm. The accession codes for all TF ChIP-seq experiments and PFM IDs are available in **Supplementary Data Sets 1a and 2b-d**.

ites. danucleotide-q data set by nsidered only hich DNase I counting the position. as a vector

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

DHSs were estimated on the basis of the DNase I raw signal. First we used the F-seq software<sup>35</sup> to create smoothed DNase-seq signals using Parzen density estimates. Then we fit the smoothed signal  $x^{\text{fseq}}$  to a gamma distribution,

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

by evaluating  $\kappa$  and  $\theta$  on the basis of mean and s.d. estimates. Finally, we found the enriched regions (DHSs) by establishing a cutoff based on a  $P$  value of 0.01 (refs. 1, 35). We refer to DHSs as a set of genomic intervals

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

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

**Estimation of DNase-seq sequence bias.** We used two approaches to estimate sequence bias in DNase-seq experiments: (1) aligned reads inside DHSs from DNase-seq experiments (termed DHS sequence bias) following He *et al.*<sup>15</sup> and (2) all aligned reads for naked DNA experiments (termed naked DNA sequence bias) following Yardimci *et al.*<sup>12</sup>. The observed cleavage score for a  $k$ -mer  $w$  corresponds to the number of DNase I cleavage sites centered at  $w$ . The background cleavage score is defined by the total number of times  $w$  occurs. Then the bias estimation is computed as the ratio between the observed and background cleavage scores. Mathematical formalizations of the bias estimation are made on the basis of the DHS sequence bias approach.

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

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

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

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

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

Finally, the cleavage bias  $b_i^s$  for a genomic position  $k + 1 \leq i \leq N - k + 1$ , given that  $w = G^s[i - (k/2) \dots i + (k/2)]$ , can be calculated as

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

where  $O^s$  indicates the total number of  $k$ -mers in DHSs,

Please make these two fractions look equal, if possible. k/2 (together) is preferable.

$$i = 1 \dots N_k$$

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

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

The bias score  $b_i^s$  represents how many times the  $k$ -mer sequence  $G^s[i - (k/2) \dots i + (k/2) + 1]$  was cleaved by the DNase I enzyme in comparison to its total occurrence in (1) DHSs (DHS sequence bias approach) and (2) the entire genome (naked DNA sequence bias approach). As observed by He *et al.*<sup>15</sup>, a 6-mer bias model captures more information than models with  $k < 6$ , and the amount of information gained in models with  $k > 6$  is not significant. Therefore, in this study, all analyses were performed using a 6-mer bias model.

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

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

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

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

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

Finally, the bias-corrected DNase-seq genomic signal ( $y$ ) was obtained by applying

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

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

**Computational footprinting methods.** In this section we present an overview of the computational footprinting methods used in

Q4. If you copy and paste the link in the browser it works. For some reason beyond my understanding, the link does not work by clicking in the pdf file. I would like to keep the link as it is. However, if you really need a clickable working link you can separate the file name from the rest of the URL, which reads: "We obtained the footprint predictions (all.footprints.gz)<ref 4> for K562 cells (DU) from ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/supplementary/integration\_data\_jan2011/byDataType/footprints/jan2011/."

the parameters from the original publication: the flanking component length varied between 3 and 10 bp, and the central footprint-region length varied between 6 and 40 bp. Afterward the footprints were filtered by a false discovery rate (FDR) of 1%, which was estimated on the basis of the FS distribution in each cell type<sup>4</sup>. Finally, we considered only predictions that occurred within DNase-seq hotspots, evaluated using the method first described by Sabo *et al.*<sup>37</sup>. We obtained all hotspots generated by Stamatoyannopoulos's lab in ENCODE<sup>1</sup> for cell types GM12878 (wgEncodeEH000492; GSM736496 and GSM736620), H1-hESC (wgEncodeEH000496; GSM736582) and K562 (wgEncodeEH000484; GSM736629 and GSM736566). We refer to this framework as the Neph framework.

**Boyle method.** Boyle *et al.*<sup>5</sup> designed a segmentation approach based on using hidden Markov models (HMMs) to predict footprints in specific DNase I cleavage patterns. Briefly, the HMM uses a normalized DNase-seq cleavage signal to find regions with depleted DNase I digestion (footprints) between two peaks of intense DNase I cleavage. Such a pattern reflects the inability of the DNase I nuclease to cleave sites where there are proteins bound. As the DNase-seq profiles required a nucleotide-resolution signal, which is usually noisy, the authors used a Savitzky-Golay smoothing filter to reduce noise and to estimate the slope of the DNase-seq signal<sup>38</sup>. Their HMM had five states, with specific states to identify the decrease or increase in DHS signals around the peak-dip-peak region. Because no source code or software was provided, we used footprint predictions from Boyle *et al.*<sup>5</sup> available at <http://fureylab.web.unc.edu/datasets/footprints/>. We refer to this method as the Boyle method.

**CENTPEDE.** CENTPEDE is a site-centric approach that gathers experimental and genomic information around MPBSs. It then uses a Bayesian mixture model approach to label each retrieved site as bound or unbound<sup>9</sup>. The experimental and genomic data used include DNase-seq, position weight matrix (PWM) bit-score, sequence conservation and distance to the nearest transcription start site. We generated the experimental data input by fetching the raw DNase-seq signal surrounding a 200-bp window centered on each MPBS. Additionally, to create the genomic data input, we obtained the PhastCons conservation score (placental mammals on the 46-way multiple alignment)<sup>39</sup> and Ensembl gene annotation from ENCODE<sup>1,40</sup> to create the prior probabilities in addition to the PWM bit-scores.

CENTPEDE software was obtained at <http://centipede.uchicago.edu/> and executed to generate posterior probabilities of regions being bound by TFs. We had previously observed that CENTPEDE is sensitive to certain parameters. Therefore, we defined CENTPEDE parameterization with an extensive computational evaluation described by Gusmao *et al.*<sup>8</sup>.

**Cuellar method.** Cuellar-Partida *et al.*<sup>10</sup> proposed a site-centric method to include DNase-seq data as priors for the detection of active TFBSs. It is based on a probabilistic classification approach to compute better log-posterior odds score than the ones observed with purely sequence-based approaches. We applied this method as described by Cuellar-Partida *et al.*<sup>10</sup>. We created a smoothed

Q5. The author moved the URL. Please update this URL to: [http://tlbailey.bitbucket.org/supplementary\\_data/Cuellar2011/](http://tlbailey.bitbucket.org/supplementary_data/Cuellar2011/)

the priors to the program FIMO<sup>41</sup> to obtain the predictions. We refer to this method as the Cuellar method.

**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. On the basis of the observation that DNase-seq cuts from the double-hit protocol are strand-specific, Wellington considers only 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 obtained Wellington's source code from <https://github.com/Wellington-Seq/Wellington>. By default, the flanking regions are 30 bp upstream and 30 bp downstream, footprint sizes varying between 6 and 40 bp with 1-bp steps, and a shoulder size (flanking regions) of 35 bp.

**Protein interaction quantification.** PIQ is a site-centric method that uses a Gaussian process to model and smooth the footprint profiles around candidate MPBSs<sup>11</sup> ( $\pm 100$  bp). Active footprints are estimated with an expectation propagation algorithm. Finally, PIQ indicates the set of motifs for which footprint signals are distinguishable from noise to reduce the set of candidate TFs. We obtained PIQ implementation from <http://piq.csail.mit.edu> and executed it with the script `python piq.py`. The DNase-seq signal was generated with the script `bam2rdata.r`. The footprint signal was created using the script `bam2rdata.r`. The footprints were detected with the script `perftf.r`.

**Footprint mixture (FLR).** Yard1mc1 *et al.*<sup>12</sup> proposed a site-centric method based on a mixture of multinomial models to detect active and 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 either initialized with 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 ( $\pm 25$  bp) as input. DNase-seq sequence bias is estimated for 6-mers on the basis of the DNA sequences extracted within the same regions in which the cuts were retrieved. Method implementation was obtained from [https://ohlerlab.mdc-berlin.de/software/FootprintMixture\\_109/](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 a TFBS (PadLen) was set to the default value of 25 bp. Also, we used the expectation maximization to re-estimate background during training (argument Fixed set to FALSE). We refer to this method as the FLR method.

**DNase2TF.** DNase2TF is a segmentation approach based on a binomial *z*-score that evaluates the depletion of DNase-seq reads around the candidate footprints<sup>7</sup>. In its 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.*<sup>8</sup> proposed the segmentation method HINT as an extension of the Boyle method<sup>5</sup>. 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 the HINT method starts with normalization of the DNase I cleavage signal using within- and between-data set normalizations. Then the slope of the normalized signals is evaluated to identify increases and decreases in the DNase-seq signal. Next an HMM is trained in a supervised manner (maximum likelihood) on the basis of a single manually annotated genomic region. To aid in such manual annotation, the normalized and slope signals are used in combination with MPBSs for all available PFMs in the repositories JASPAR<sup>32</sup> and UniPROBE<sup>33</sup>. Finally, the Viterbi algorithm is applied to the trained HMMs inside regions consisting of DHSs extended by 5,000 bp upstream and downstream. All parameters were set as described by Gusmao *et al.*<sup>8</sup>.

We made two modifications to the method described by Gusmao *et al.*<sup>8</sup>. First, to perform a standardized comparison, we modified HINT to allow only DNase-seq data. The modified HMM model contained 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 before normalization steps. These modifications required retraining of the HMMs. For this, we used the same manual annotation described by Gusmao *et al.*<sup>8</sup>. The novel methods and trained models are available as a command-line tool at <http://costalab.org/publications-2/hint-bc/>.

**BinDNase.** BinDNase is a site-centric method based on logistic regression that is used to predict active and inactive MPBSs<sup>13</sup>. The algorithm starts with the base pair-resolution DNase-seq signal around MPBSs ( $\pm 100$  bp) 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 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 that requires TF ChIP-seq examples for training. We also point out that BinDNase did not execute successfully for 19 TFs in our evaluation data set (POU5F1, REST, RFX5, SP1, SP2, SRF, TCF12 and ZNF143 binding in H1-hESC cells; ARID3A, CTCF, IRF1, MEF2A, PU1, REST, RFX5, SP1, SP2, STAT2 and ZNF263 binding in K562 cells) given the lack of ChIP-seq data. The BinDNase segmentation scheme termed the footprint score (FS), which is based on a scoring metric from the footprinting methodology proposed by Neph *et al.*<sup>4</sup>. The FS statistic is defined as

$$FS_{MPBS_i} = - \frac{\log \left( \frac{n_{C,i}}{n_{R,i} + 1} + \frac{n_{G,i}}{n_{L,i} + 1} \right)}{\log 2}$$

where  $MPBS_i = [m_i, n_i]$  is the *i*th MPBS, which extends from genomic positions  $m_i$  to  $n_i$  and  $MPBS_i = (m+n)/2$ . 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

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

**TC-Rank.** The site-centric method referred to here as the tag count (TC) corresponds to the number of DNase I cleavage hits in a 200-bp window around predicted TFBSs as defined by He *et al.*<sup>15</sup>. This can be expressed as

$$TC_{MPBS_i} = \sum_{j=MPBS_i-100}^{MPBS_i+99} x_j$$

Both TC and FS can be used as quality scores for footprints. However, their respective methods (TC-Rank and FS-Rank) consist of attributing these quality scores to each MPBS and evaluating the performance at these ranked MPBSs. This observation also holds for the PWM-Rank method described below.

**Evaluation.** *MPBSs.* Method evaluation was performed with site-centric binding site statistics. For this we generated PWMs from PFMs by evaluating the information content of each position and performing background nucleotide frequency correction<sup>42</sup>. This was done using Biopython<sup>43</sup>. Then we created MPBSs by matching all PWMs against the human (hg19) and mouse (mm9) genomes using the fast-performance motif-matching tool MOODS<sup>44</sup>. This procedure produces PWM bit-scores for every match. We determined a bit-score cutoff threshold by applying the dynamic programming approach described by Wilczynski *et al.*<sup>45</sup> with an 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 in a baseline method referred to here as PWM-Rank.

**Method comparison.** Methods were evaluated using a site-centric approach<sup>10</sup> that combined 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, and MPBSs without ChIP-seq evidence are considered false TFBSs. Every TF prediction that overlaps a true TFBS is considered a correct prediction (true positive), and every prediction that overlaps a false TFBS is considered an incorrect prediction (false positive). Therefore, true negatives and false negatives 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 the specificity of the method is decreased. The AUC metric was evaluated at 100%, 10% and 1% FPR. We also evaluated the AUPR. This metric is indicated for problems with imbalanced data sets (distinct numbers of positive and negative examples)<sup>20,46</sup>.

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 and then ranked all nonpredicted sites by their tag count. To present a fair comparison, we also applied this approach 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 and 0.99 for the detection of footprints in all site-centric methods. We performed additional experiments to select the best threshold for each method (**Supplementary Fig. 8**).

Our TF ChIP-seq-based comparative experiments comprised the following three evaluation scenarios (He, benchmarking and comprehensive data sets). All evaluation statistics and details on method performance are available in **Supplementary Data Set 1**.

**He data set.** To replicate the analysis performed by He *et al.*<sup>15</sup>, we analyzed DNase-seq data from cell types K562 (UW), LNCaP (DU) and m3134 (UW) for 36 TFs and evaluated the methods PWM, FS, TC, HINT, HINT-BC and HINT-BCN.

**Benchmarking data set.** For comparative analysis of several competing methods, we selected the two cell types with the highest numbers of TF ChIP-seq data sets evaluated in our study: K562 (DU), with 59 TFs, and H1hESC (DU), with 29 TFs. We were therefore able to make use of predictions provided by Gusmao *et al.*<sup>8</sup> and Boyle *et al.*<sup>5</sup>, including evaluation of the PWM, Boyle, Cuellar, CENTIPEDE, HINT and Neph methods. For this data set, we estimated novel footprints for the FS, TC, HINT-BC, HINT-BCN, DNase2TF, PIQ, Wellington and FLR methods, which were not previously evaluated.

**Comprehensive data set.** Lastly, we compiled a comprehensive data set containing 233 cell types, including cellular background of 144 TF ChIP-seq and 15 DNase-seq data sets. These data were 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 by Yardimci *et al.*<sup>12</sup>, ChIP-seq evaluation of putative TFBSs may present biases regarding the fact that ChIP-seq data alone are not able to distinguish direct from indirect binding events. Consequently, we performed an evaluation procedure that combined MPBSs with differentially expressed genes from two cell types. The method evaluated the association of the quality of footprints overlapping particular motifs and the expression of the relevant TF.

We used limma<sup>47</sup> to perform between-array normalization on expression of H1-hESC, K562 and GM12878 cells and obtain fold-change estimates. Then we retrieved all nonredundant PFMs from JASPAR in which the gene symbol was a perfect match with genes present in the array platform. This led us to 143 PFMs (**Supplementary Data Set 2b–d**). We applied genome-wide motif matching using these PFMs.

Next we evaluated the FLR<sup>12</sup> score, TC<sup>15</sup> and FS<sup>15</sup> for the footprints of each evaluated method that intersected with MPBSs of a particular motif. We considered only the footprints within DHSs that were common between the two cell types being evaluated, as described by Yardimci *et al.*<sup>12</sup>. We expected that TFs expressed in cell type A would present higher values for 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

Change to "Yardimci" to match the main text.



versa. We used a two-sample Kolmogorov-Smirnov (KS) test to assess the difference between each metric's 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 greater differences. As the KS score does not indicate the direction of the change in distribution, we obtained a signed version by multiplying the KS statistic by  $-1$  in cases where the median of A was less than the median of B. We calculated the Spearman correlation between the signed KS test statistic and the expression fold change for each TF (Supplementary Figs. 1 and 2). Positive values indicated an association between expression of TFs and quality of footprint predictions. We call this correlation FLR-Exp. Results for FLR-Exp analysis are summarized in Supplementary Data Set 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 or TF. As previously indicated by Sung *et al.*<sup>7</sup>, TFs with shorter binding times characteristically have fewer DNase-seq cuts (protection) surrounding the binding site. More formally, the protection score for a set of MPBSs is defined as

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

where  $\text{MPBS} = \{\text{MPBS}_1, \dots, \text{MPBS}_N\}$  is a set of binding sites for a given motif;  $\text{MPBS}_i = [m_i, n_i]$  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).

In short, the protection score indicates the average difference between DNase-seq counts in the flanking region and DNase-seq counts within the MPBS. Positive values indicate protection in the flanking regions, whereas values close to zero and negative values indicate no protection. The protection score is similar to the FS<sup>15</sup>. The main difference is that the FS measures the ratio between reads in flanking and binding sites, whereas the protection score measures the difference. Finally, because we were interested in using the protection score as a measure of quality for a given TF and set of footprint predictions, we evaluated only MPBSs overlapping footprints for a given cell type. The DNase-seq count values had been previously corrected for DHS sequence bias and coverage differences. Results for protection scores are provided in Supplementary Data Set 1.

**Statistical methods.** We used the nonparametric Friedman-Nemenyi hypothesis test<sup>48</sup> to compare the AUC and AUPR of the methods regarding all data set combinations (TFs versus cell types). Such a test provides a rank of the methods as well as the statistical significance of the outperformance of a particular method. All correlations are based on Spearman values. All reported  $P$  values were corrected using the Benjamini-Hochberg method<sup>49</sup>.

**Code availability.** Software, custom code, benchmarking data, DNase-seq sequence bias estimates and additional graphical results are available at <http://costalab.org/publications-2/hint-bc/>. HINT, HINT-BC and HINT-BCN software can be accessed directly through the Regulatory Genomics Toolbox website at <http://www.regulatory-genomics.org/hint/>.

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. & Bullyk, 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. & Furey, 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 Bioinformatics* **10**, 82 (2009).
46. Fawcett, T. An introduction to ROC analysis. *Pattern Recognit. 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).
49. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Series B Stat. Methodol.* **57**, 289–300 (1995).

# QUERY FORM

Nature Methods	
Manuscript ID	[Art. Id: 3772]
Author	
Editor	
Publisher	

## AUTHOR:

The following queries have arisen during the editing of your manuscript. Please answer queries by making the requisite corrections directly on the galley proof. It is also imperative that you include a typewritten list of all corrections and comments, as handwritten corrections sometimes cannot be read or are easily missed. Please verify receipt of proofs via e-mail

Query No.	Nature of Query
Q1	Please carefully check the spelling and numbering of all author names and affiliations.
Q2	Please check that all funders have been appropriately acknowledged and that all grant numbers are correct.
Q3	Please check that the Competing Financial Interests declaration is correct as stated. If you declare competing interests, please check the full text of the declaration (at the end of the main references section) for accuracy and completeness.
Q4	Please check. link not working
Q5	Url not working Please check.
Q6	Labels are overlapping. Please check.
	<div>All queries were answered directly in the manuscript and are numbered accordingly. Furthermore, additional comments and corrections were made directly in the manuscript.</div>