

1 **Beyond accessibility: ATAC-seq footprinting unravels kinetics**  
2 **of transcription factor binding during zygotic genome**  
3 **activation**

4 **Authors**

5 Mette Bentsen<sup>1</sup>, Philipp Goymann<sup>1</sup>, Hendrik Schultheis<sup>1</sup>, Kathrin Klee<sup>1</sup>, Anastasiia Petrova<sup>1</sup>,  
6 René Wiegandt<sup>1</sup>, Annika Fust<sup>1</sup>, Jens Preussner<sup>1,3</sup>, Carsten Kuenne<sup>1</sup>, Thomas Braun<sup>2,3</sup>,  
7 Johnny Kim<sup>2,3</sup>, Mario Looso<sup>1,3</sup>

8 **Affiliation**

9 <sup>1</sup> Bioinformatics Core Unit (BCU), Max Planck Institute for Heart and Lung Research, Bad  
10 Nauheim, Germany

11 <sup>2</sup> Department of Cardiac Development and Remodeling, Max-Planck-Institute for Heart and  
12 Lung Research, Bad Nauheim, Germany

13 <sup>3</sup> German Centre for Cardiovascular Research (DZHK), Partner site Rhein-Main, Frankfurt am  
14 Main, 60596 Germany

15 **Corresponding author email address**

16 [mario.looso@mpi-bn.mpg.de](mailto:mario.looso@mpi-bn.mpg.de)

## 17 **Abstract**

18 While footprinting analysis of ATAC-seq data can theoretically enable investigation of  
19 transcription factor (TF) binding, the lack of a computational tool able to conduct different  
20 levels of footprinting analysis has so-far hindered the widespread application of this method.  
21 Here we present TOBIAS, a comprehensive, accurate, and fast footprinting framework  
22 enabling genome-wide investigation of TF binding dynamics for hundreds of TFs  
23 simultaneously. As a proof-of-concept, we illustrate how TOBIAS can unveil complex TF  
24 dynamics during zygotic genome activation (ZGA) in both humans and mice, and explore how  
25 zygotic Dux activates cascades of TFs, binds to repeat elements and induces expression of  
26 novel genetic elements. TOBIAS is freely available at: <https://github.com/loosolab/TOBIAS>.

## 27 **Keywords**

28 Footprinting, ATAC-seq, epigenetics, transcription factors, ZGA, Dux

## 29 **Background**

30 Epigenetic mechanisms governing chromatin organization and transcription factor (TF)  
31 binding are critical components of transcriptional regulation and cellular transitions. In recent  
32 years, rapid improvements of pioneering sequencing methods such as ATAC-seq (Assay of  
33 Transposase Accessible Chromatin) <sup>1</sup>, have allowed for systematic, global scale investigation  
34 of epigenetic mechanisms controlling gene expression. While ATAC-seq can uncover  
35 accessible regions where TFs might bind, reliable identification of specific TF binding sites  
36 (TFBS) still relies on chromatin immunoprecipitation methods such as ChIP-seq. However,  
37 ChIP-seq methods require high input cell numbers, are limited to one TF per assay, and are  
38 further restricted to TFs for which antibodies are readily available. Latest improvements of  
39 ChIP based methods <sup>2</sup> can circumvent some of these technical drawbacks, but the limitation  
40 of only being able to identify binding sites of one TF per assay persists. Therefore, it remains  
41 costly, or even impossible, to study the binding of multiple TFs in parallel.

42 Current limits to the investigation of TF binding become particularly apparent when  
43 investigating processes involving a very limited number of cells such as preimplantation  
44 development (PD) of early zygotes. During PD, the fertilized egg forms the zygote, which  
45 undergoes a series of cell divisions to finally constitute the blastocyst, a structure built by the  
46 inner cell mass (ICM) and trophectoderm (Figure 1a). Within this process, maternal and  
47 paternal mRNAs are degraded prior to zygotic genome activation (ZGA) (reviewed in <sup>3</sup>), a  
48 transformation which eventually leads to the transcription of thousands of genes <sup>4</sup>. Integration  
49 of multiple omics-based profiling methods have revealed a set of key TFs that are expressed  
50 at the onset of and during ZGA including Dux <sup>5, 6</sup>, Zscan4 <sup>7</sup>, and other homeobox-containing  
51 TFs <sup>8</sup>. However, due to the limitations of ChIP-seq, the exact genetic elements bound and  
52 regulated by different TFs during PD remain to be fully discovered. Consequently, the global  
53 network of TF binding dynamics throughout PD remains mostly obscure.

54 A computational method known as *digital genomic footprinting* (DGF) <sup>9</sup> has emerged as an  
55 alternative means, which can overcome some of the limits of ChIP-based methods. DGF is a  
56 computational analysis of chromatin accessibility assays such as ATAC-seq, which employs  
57 DNA effector enzymes that only cut accessible DNA regions. Similarly to nucleosomes, bound  
58 TFs hinder cleavage of DNA, resulting in defined regions of decreased signal strength within  
59 larger regions of high signal - known as *footprints* <sup>10</sup> (Figure 1b).

60 Surprisingly, although this concept shows considerable potential to survey genome-wide  
61 binding of multiple TFs in parallel from a single experiment, DGF analysis is rarely applied  
62 when investigating TF binding mechanisms. The skepticism towards DGF has been driven by  
63 the discovery that enzymes used in chromatin accessibility assays (e.g. DNase-I) are biased  
64 towards certain sequence compositions, an effect which has been well characterized for  
65 DNase-seq <sup>11, 12</sup>. The influence of Tn5 transposase bias in the context of ATAC-seq  
66 footprinting has, however, only been described very recently <sup>13, 14</sup> and still represents an  
67 uncertainty during discovery of true footprints. Besides the identification of footprints,  
68 comparing footprints across biological conditions remains challenging as well. While there  
69 have been efforts to estimate differential TF binding on a genomewide scale <sup>15, 16</sup>, investigation  
70 of epigenetic processes often require more in-depth information on the individual differentially  
71 bound TFBS and genes targeted by these TFs, which is not provided by these methods.  
72 Furthermore, many footprinting methods suffer from performance issues due to missing  
73 support for multiprocessing, inflexible software architecture prone to software dependency  
74 issues, and the use of non-standard file-formats. These obstacles complicate the assembly of  
75 different tools for advanced analysis workflows. Consequently, despite its compelling  
76 potential, these issues have rendered footprinting on ATAC-seq cumbersome to apply to  
77 biological questions. Essentially, a comprehensive framework enabling large-scale ATAC-seq  
78 footprinting is missing.

79 Here, we describe and apply TOBIAS (**T**ranscription factor **O**ccupancy prediction **B**y  
80 **I**nvestigation of **A**TAC-seq **S**ignal), a comprehensive computational framework that we

81 created for footprinting analysis (Figure 1c). TOBIAS is a collection of command-line tools,  
82 which provides functionality to perform all levels of footprinting analysis including bias  
83 correction, footprinting, and comparison between conditions (Supp. Figure 1; Footprinting  
84 pipeline). Furthermore, TOBIAS includes a variety of auxiliary tools such as TF network  
85 inference and visualization of footprints, which can be combined for more targeted  
86 downstream analysis (Supp. Figure 1; Supporting tools). To couple individual tools, we provide  
87 scalable analysis workflows implemented in Snakemake <sup>17</sup> and NextFlow <sup>18</sup>, including a cloud  
88 computing compatible version making use of the de.NBI cloud <sup>19</sup>. These pipelines utilize a  
89 minimal input of ATAC-seq reads, TF motifs and genome information to enable complete  
90 footprinting analysis and comparison of TF binding even for complex experimental designs  
91 (e.g. time series).

## 92 **Results**

### 93 **Classification and validation of TOBIAS**

94 A computational DGF framework able to perform footprinting on ATAC-seq data and handle  
95 complex experimental designs autonomously does not exist. Nonetheless, to demonstrate the  
96 advantages of TOBIAS, we compared the individual framework features to published  
97 footprinting tools for ATAC-seq footprinting where applicable. While we found that some  
98 functionalities are overlapping between tools, we found a substantial set of features  
99 exclusively covered by TOBIAS (Table 1). As sequencing costs will continue to decrease,  
100 allowing for ever more data to be created, it is worth noting that TOBIAS is the only tool  
101 supporting differential footprinting for more than two conditions. Additionally, TOBIAS is one  
102 of only two footprinting tools applying multiprocessing to speed up computation, resulting in  
103 the lowest runtime among the compared set of tools.

104 To compare the footprinting capabilities of individual tools, we utilized 218 paired ChIP-seq /  
105 ATAC-seq datasets across four different cell types. Here, the ChIP-seq data represents the

106 true binding sites for each TF, which we used for validating each method after application to  
107 the matched ATAC-seq data (see Methods; Validation). In terms of Tn5 bias correction, as  
108 well as visualization of aggregate footprints, we found that TOBIAS clearly outperforms other  
109 bias-correction tools in uncovering footprints and thereby distinguishing between  
110 bound/unbound sites (Supp. Figure 2a, Supp. File 1). For the task of protein binding prediction,  
111 we found that TOBIAS significantly outperformed the other de novo tools HINT-ATAC, PIQ  
112 and Wellington (Supp. Figure 2b) and performed equally well as msCentipede overall (Supp.  
113 Figure 2c). Notably, TOBIAS also showed robust performance across individual cell types  
114 (Suppl. Figure 2d). Looking at individual TFs, TOBIAS outperforms msCentipede for factors  
115 with a notable gain of footprints after Tn5 bias correction (Supp. Figure 2e), once again  
116 highlighting the advantage of taking Tn5 bias into account. Although msCentipede implements  
117 a motif centric learning approach, which can take TF specific binding patterns into account, it  
118 did not yield overall higher accuracy in comparison to TOBIAS. Additionally, the approach of  
119 building individual TF models took 300 times longer to compute than performing footprinting  
120 using TOBIAS (Supp. Figure 2f and Table 1). Such learning approaches are therefore greatly  
121 limited in the number of TFs and conditions that can realistically be included in an analysis. In  
122 conclusion, we find that the TOBIAS framework shows unprecedented accuracy and  
123 performance in the field of ATAC-seq footprinting.

124 In order to confirm the improvement of footprint detection after Tn5 bias correction, we made  
125 use of another exemplary dataset derived from hESC<sup>20</sup>. Importantly, besides cases where the  
126 footprint was hidden by Tn5 bias (Supp. Figure 3a; ZSCAN4), we also identified TFs for which  
127 the motif itself disfavors Tn5 integration, thereby creating a false-positive footprint in  
128 uncorrected signals, which disappears after Tn5 bias correction (Supp. Figure 3a; HLX).  
129 Utilizing a footprint depth metric as described by<sup>16</sup> (Supp. Figure 3b) across uncorrected,  
130 expected and corrected Tn5 signals, we found a high correlation between uncorrected and  
131 expected footprinting depths (Supp. Figure 3c). In contrast, this effect vanished after TOBIAS  
132 correction (Supp. Figure 3d), effectively uncovering TF footprints which were superimposed

133 by Tn5 bias. From a global perspective, taking 590 TFs into account, TOBIAS generated a  
134 measurable footprint for 64% of the TFs (Supp. Figure 3e). This is in contrast to previous  
135 reports wherein it has been suggested that only 20% of all TFs leave measurable footprints  
136 <sup>16</sup>. To summarize, we found that TOBIAS exceeded other tools in terms of uncovering  
137 footprints and correctly identifying bound TF binding sites.

### 138 **Footprinting uncovers transcription factor binding dynamics in mammalian ZGA**

139 To demonstrate the full potential of TOBIAS, in particular in the investigation of processes  
140 involving only few cells, we analyzed a series of ATAC-seq datasets derived from both human  
141 and murine preimplantation embryos at different developmental stages ranging from 2C, 4C,  
142 8C to ICM in addition to embryonic stem cells of their respective species <sup>20, 21</sup>. Altogether,  
143 TOBIAS was used to calculate footprint scores for a list of 590 and 464 individual TFs across  
144 the entire process of PD of human and mouse embryos, respectively. After clustering TFs into  
145 co-active groups within one or multiple developmental timepoints, we first asked whether the  
146 predicted timing of TF activation reflects known processes in human PD. Intriguingly, we found  
147 10 defined clusters of specific binding patterns, the majority of which peaked between 4C and  
148 8C, fully concordant with the transcriptional burst and termination of ZGA (Figure 2a and Supp.  
149 Table 1).

150 Two clusters of TFs (Cluster 1+2; n=83) displayed highest activity at the 2-4C stage and  
151 strongly decreased thereafter, suggesting that factors within these clusters are likely involved  
152 in ZGA initiation. We set out to classify these TFs, and observed a high overlap with known  
153 maternally transferred transcripts <sup>22</sup> (LHX8, BACH1, EBF1, LHX2, EMX1, MIXL1, HIC2,  
154 FIGLA, SALL4, ZNF449), explaining their activity before ZGA onset. Importantly, DUX4 and  
155 DUXA, which are amongst the earliest expressed genes during ZGA <sup>5, 6</sup>, were also contained  
156 in these clusters. Additional TFs included HOXD1, which is known to be expressed in human  
157 unfertilized oocytes and preimplantation embryos <sup>23</sup> and ZBTB17, a TF mandatory to generate  
158 viable embryos <sup>24</sup>. Cluster 6 (n=67) displayed a particularly prominent 8C specific signature,  
159 that harbored well known TFs involved in lineage specification such as PITX1, PITX3, SOX8,

160 MEF2A, MEF2D, OTX2, PAX5 and NKX3.2. Furthermore, overlapping TFs within Cluster 6  
161 with RNA expression datasets ranging from the germinal vesicle to cleavage stage <sup>5</sup>, 12  
162 additional TFs (FOXJ3, HNF1A, ARID5A, RARB, HOXD8, TBP, ZFP28, ARID3B, ZNF136,  
163 IRF6, ARGFX, MYC, ZSCAN4) were confirmed to be exclusively expressed within this time  
164 frame. Taken together, these data show that TOBIAS reliably uncovers massively parallel TF  
165 binding dynamics at specific time points during early embryonic development.

### 166 **Transcription factor scores correlate with footprints and gene expression**

167 To confirm that TOBIAS-based footprinting scores are indeed associated with leaving *bona*  
168 *fade* footprints we utilized the ability to visualize aggregated footprint plots as implemented  
169 within the framework. Indeed, bias corrected footprint scores were highly congruent with  
170 explicitly defined footprints (Figure 2b) of prime ZGA regulators at developmental stages in  
171 which these have been shown to be active <sup>7</sup>. For example, footprints associated with DUX4,  
172 a master inducer of ZGA, were clearly visible from 2C-4C, decreased from 8C onwards and  
173 were completely lost in later stages, consistent with known expression levels <sup>20</sup> and ZGA onset  
174 in humans. Footprints for ZSCAN4, a primary DUX4 target <sup>5</sup>, were exclusively visible at the  
175 8C stage. Interestingly, GATA2 footprints were exhibited from 8C to ICM stages which is in  
176 line with its known function in regulating trophoblast differentiation <sup>25</sup>. As expected, CTCF  
177 creates footprints across all timepoints. Strikingly, we observed that these defined footprints  
178 were not detectable without TOBIAS mediated Tn5 bias correction (Supp. Figure 3f). These  
179 data show that footprint scores can be reliably confirmed by footprint visualizations, which  
180 further allow to infer TF binding dynamics.

181 To test if the global footprinting scores of individual TFs correlate with the incidence and level  
182 of their RNA expression, we matched them to RNA expression datasets derived from  
183 individual timepoints throughout zygotic development, taking TF motif similarity into account.  
184 Indeed, we found that TOBIAS scores for the majority of TFs either correlated well with the  
185 timing of their expression profiles or displayed a slightly delayed activity after expression



186 peaked (Supp. Figure 4a). This is important because it shows that in conjunction with  
187 expression data, TOBIAS can unravel the kinetics between TF expression (mRNA) and the  
188 actual binding activity of their translated proteins. The value of this added information becomes  
189 particularly apparent when analyzing activities of TFs that did not correlate with the timing of  
190 their RNA expression (Supp. Figure 4a; not correlated).

191 For example, within the non-correlated cluster 13 TFs were identified which are of putative  
192 maternal origin <sup>22</sup> including SALL4. In mice, Sall4 protein is maternally contributed to the  
193 zygote, subsequently degraded at 2C and then reexpressed after zygotic transcription has  
194 initiated <sup>26</sup>. Consistent with this, SALL4 expression increases dramatically from 8C onwards  
195 (Supp. Table 2). In contrast to the expression values, TOBIAS predicted SALL4 to have the  
196 highest activity in 2C and second-highest activity in hESC (on-off-on-pattern), which is in line  
197 with the presence of maternal SALL4 in the zygote. These data show that TOBIAS can predict  
198 true on-off-on-patterns, and can infer significant insight into TF activities, in particular for those  
199 where determining their expression patterns alone does not suffice to explain when they exert  
200 their biological function.

## 201 **Differential footprint analysis reveals functional divergence between human and mouse**

### 202 **ZGA**

203 The timing of ZGA varies between mice (2C) and humans (4C to 8C) (reviewed in <sup>27 28</sup>). By  
204 integrating the TOBIAS scores from human and mouse (Supp. Figure 4b and Supp. Table 3),  
205 and instrumentalizing the capability of TOBIAS to generate differential TF binding plots for all  
206 time points automatically, we investigated similarities and differences of PD between these  
207 species. Firstly, reflecting the shift of ZGA onset, we identified 30 TFs which appeared to be  
208 ZGA specific in both human and mouse (Figure 2c) including several homeobox factors which  
209 already have described functions within ZGA <sup>29</sup> as well as ARID3A which has been shown to  
210 play a role in cell fate decisions in creating trophectoderm <sup>30</sup>.

211 Next, we used the differential TF binding plots to display differences in ZGA at the transition  
212 between 2C and 4C in mouse (Figure 3a), and human 8C and ICM (Figure 3b) (Supp. File 2  
213 + 3 for all pairwise comparisons). In mice, we observed a shift of Obox-factor activity in 2C to  
214 an activation of Tead (Tead1-4) and AP-2 (Tfap2a/c/e) motifs in 4C. Notably, AP-2/Tfap2c is  
215 required for normal embryogenesis in mice <sup>31</sup> and was also recently shown to act as a  
216 chromatin modifier that opens enhancers proximal to pluripotency factors in human <sup>32</sup>. We  
217 observed a similar shift of TF activity for homeobox factors such as PITX1-3, RHOXF1, CRX  
218 and DMBX1 at the human 8C stage towards higher scores in ICM for known pluripotency  
219 factors such as POU5F1 (OCT4) and other POU-factors. Taken together, these results  
220 highlight the ability of TOBIAS to capture differentially bound TFs, not only across the whole  
221 timeline, but also between individual conditions and species.

222 Throughout the pairwise comparisons, we observed that TFs from the same families often  
223 display similar binding kinetics within species, which is not surprising since they often possess  
224 highly similar binding motifs (Figure 3a right). To characterize TF similarity, TOBIAS clusters  
225 TFs based on the overlap of TFBS within investigated samples (Figure 3c+3d). This enables  
226 quantification of the similarity and clustering of individual TFs that appear to be active at the  
227 same time. Thereby, we observed a group of homeobox motifs which cluster together with  
228 more than 50% overlap of their respective binding sites in mouse (Figure 3c). In contrast, other  
229 TFs such as Tead and AP-2 cluster separately, indicating that these factors utilize independent  
230 motifs (Supp. File 2+3). While this might appear trivial, this clustering of TFs in fact also  
231 highlights differences in motif usage between human and mouse. One prominent example is  
232 the RHOXF1 motif, which shows high binding-site overlap with Obox 1/3/5 and Otx2 binding  
233 sites in mouse (Figure 3c; ~60% overlap), but does not cluster with OTX2 in human (Figure  
234 3d; ~35% overlap). This observation suggests important functional differences of RHOX/Rhox  
235 TFs between mice and humans. In support of this hypothesis RHOXF1, RHOXF2 and  
236 RHOXF2B genes are exclusively expressed at 8C and ICM in humans, whereas Rhox factors  
237 are not expressed in corresponding developmental stages of preimplantation in mouse (Supp.

238 Table 4). Conceivably, this observation, together with the finding that murine Obox factors  
239 share the same motif as RHOX-factors in humans, suggests that Obox TFs might function  
240 similarly to RHOX-factors during ZGA. Altogether, the TOBIAS mediated TF clustering based  
241 on TFBS overlap allows for quantification of target-similarity and divergence of TF function  
242 between motif families.

### 243 **Dux expression induces massive changes of chromatin accessibility, transcription and** 244 **TF networks**

245 Throughout the investigations of human and mouse development we became particularly  
246 interested in the Dux/DUX4 TF, which TOBIAS predicted to be one of the earliest factors to  
247 be active in both organisms (Figure 2a, Supp. Figure 4b and Supp. Table 1+3). Interestingly,  
248 despite the fact that Dux has already been proved to play a prominent role in ZGA<sup>5-7, 33, 34</sup>,  
249 there is still a poor understanding of how Dux regulates its primary downstream targets, and  
250 consequently its secondary targets, during this process. We therefore applied TOBIAS to  
251 identify Dux binding sites utilizing an ATAC-seq dataset of Dux overexpression (DuxOE) in  
252 mESC<sup>5</sup>.

253 As expected, the differential TF activity predicted by TOBIAS showed an increase in activity  
254 of Dux, Obox and other homeobox-TFs (Figure 4a, Supp. File 4). Interestingly, this was  
255 accompanied by a massive loss of TF binding for pluripotency markers such as Nanog, Pou5f1  
256 (OCT4) and Sox2 upon DuxOE, indicating that Dux renders previously accessible chromatin  
257 sites associated with pluripotency inaccessible.

258 Consistently, Dux footprints (Figure 4b; left) were clearly evident upon DuxOE. In comparison  
259 to existing bias-correction methods, we found TOBIAS to be superior in uncovering this  
260 footprint between Control and DuxOE conditions (Supp. Figure 5a). Importantly, TOBIAS  
261 additionally discriminated ~30% of all potential binding sites within open chromatin regions to  
262 be bound in the DuxOE condition, which further demonstrates the specificity of this method  
263 (Figure 4b; right). To rank the biological relevance of the individually changed binding sites

264 between control and DuxOE conditions, we linked all annotated gene loci to RNA expression.  
265 A striking correlation between the gain-of-footprint and gain-of-expression of corresponding  
266 loci was clearly observed and mirrored by the TOBIAS predicted bound/unbound state (Figure  
267 4c). Amongst the genes within the list of bound Dux binding sites (Supp. Table 5 for full Dux  
268 target list) were well known Dux targets including *Zscan4c* and *Pramef25*<sup>35</sup>, for which local  
269 footprints for Dux were clearly visible (Figure 4d). The high resolution of footprints is  
270 particularly pronounced for *Tdpoz1* which harbors two potential Dux binding sites of which one  
271 is clearly footprinted in the score track, while the other is predicted to be unoccupied (Figure  
272 4d; bottom). In line with this, *Tdpoz1* expression is significantly upregulated upon DuxOE as  
273 revealed by RNA-seq (log2FC: 6,95). Consistently, *Tdpoz1* expression levels are highest at  
274 2C in zygotes and decrease thereafter, strongly indicating that *Tdpoz1* is likely a direct target  
275 of Dux during PD both *in vitro* and *in vivo*<sup>21, 36</sup> (Supp. Table 5). Footprinting scores also  
276 directly correlated with ChIP-seq peaks for Dux in the *Tdpoz1* promoter (Supp. Figure 5b), an  
277 observation which we also found at many other positions (Examples shown in Supp. Figure  
278 5c+d).

279 Many of the TOBIAS-predicted Dux targets encode TFs themselves. Therefore, we applied  
280 the TOBIAS network module to subset and match all activated binding sites to TF target genes  
281 with the aim of inferring how these TF activities might connect. Thereby, we could model an  
282 intriguing pseudo timed TF activation network. This directed network uncovered a TF  
283 activation cascade initiated by Dux, resulting in the activation of 7 primary TFs which appear  
284 to subsequently activate 32 further TFs (first three layers depicted in Figure 4e). As Dux is a  
285 regulator of ZGA, we asked how the *in vitro* activated Dux network compared to gene  
286 expression throughout PD *in vivo*. Strikingly, the *in vivo* RNA-seq data of the developmental  
287 stages<sup>21</sup> confirmed an early 2C specific expression for Dux, followed by a slightly shifted  
288 activation pattern for all direct Dux targets except for *Rxrg* (Figure 4f). However, it is of note  
289 that *Rxrg* is significantly upregulated in the *in vitro* DuxOE from which the network is inferred  
290 (Supp. Table 5), pointing to both the similarities and differences between the *in vivo* 2C and

291 *in vitro* 2C-like stages induced by Dux. In conclusion, these data show that beyond identifying  
292 specific target genes of individual TFs, TOBIAS can infer biological insight by predicting entire  
293 TF activation networks.

294 Notably, many of the predicted Dux binding sites (40%) are not annotated to genes (Figure  
295 4g), raising the question what role these sites play in ZGA. Dux is known to induce expression  
296 of repeat regions such as LTRs<sup>5</sup> and consistently, we found that more than half of the DUX-  
297 bound sites without annotation to genes are indeed located within known LTR sequences  
298 (Figure 4g) which were transcribed both *in vitro* and *in vivo* (Figure 4h). Interestingly, we  
299 additionally found that 28% of all non-annotated Dux binding sites overlap with genomic loci  
300 encoding LINE1 elements. Although LINE1 expression does not appear to be altered in mESC  
301 cells, there is a striking pattern of increasing LINE1 transcription from 4C-8C (Figure 4h) *in*  
302 *vivo*, pointing to a possible role of LINE1 regulation throughout PD. Finally, we found a portion  
303 of the Dux binding sites which do not overlap with any annotated gene nor with putative  
304 regulatory repeat sequences, even though transcription clearly occurs at these sites (Figure  
305 4h; bottom). One example is a predicted Dux binding site on chromosome 13, which coincides  
306 with a spliced region of increased expression between control mESC/DuxOE and comparable  
307 high expression in 2C, 4C and 8C (Supp. Figure 6). These data clearly indicate the existence  
308 of novel transcribed genetic elements, the function of which remains unknown, but which are  
309 likely controlled by Dux and could play a role during PD.

310 In conclusion, TOBIAS predicted the exact locations of Dux binding in promoters of target  
311 genes, and could unveil how Dux initiates TF-activation networks and induces expression of  
312 repeat regions. Importantly, these data further show that TOBIAS can identify any TFBS with  
313 increased binding, not only those limited to annotated genes, which aids in uncovering novel  
314 regulatory genetic elements.

## 315 **Discussion**

### 316 **Footprint scores reveal true characteristics of protein binding**

317 To the best of our knowledge, this is the first application of a DGF approach to visualize gain  
318 and loss of individual TF footprints in the context of time series, TF overexpression, and TF-  
319 DNA binding for a wide-range of TFs in parallel. Importantly, we found that these advances  
320 could in large part be attributed to the framework approach we took in developing TOBIAS,  
321 which enabled us to simultaneously compare global TF binding across samples and quantify  
322 changes in TF binding at specific loci. The modularity of the framework also allowed us to  
323 apply a multitude of downstream analysis tools to easily visualize footprints and gain even  
324 more information about TF binding dynamics as exemplified by the discovery of the Dux TF-  
325 activation network.

326 The power of this framework to handle time-series data becomes especially apparent when  
327 correlating the TOBIAS-based prediction of TF binding to RNA-seq data from the same time  
328 points. For instance, TOBIAS could infer when the maternally transferred TF SALL4 is truly  
329 active while its gene expression pattern alone does not allow to make such conclusions. Along  
330 this line, TOBIAS is also powerful in circumstances where gene expression of a particular TF  
331 appears to be anticorrelated with its binding activity. It is tempting to speculate that TFs for  
332 which footprinting scores are low, even though their RNA expression is high, might act as  
333 transcriptional repressors, because footprinting relies on the premise that TFs will increase  
334 chromatin accessibility around the binding site. In support of this hypothesis, recent  
335 investigations have suggested that repressors display a decreased footprinting effect in  
336 comparison to activators<sup>37</sup>. Therefore, the integration of ATAC-seq footprinting and RNA-seq  
337 is an important step in revealing additional information such as classification TFs into  
338 repressors and activators, as well as the kinetics between expression and binding.

### 339 **Species-specific TFs use common ZGA motifs in mice and human**

340 By integration of human and murine TF activities using both differential footprinting and  
341 species-specific TFBS overlaps, our analyses revealed that the majority of TF motifs are active  
342 at corresponding timepoints of human and mouse ZGA. This is not necessarily surprising since  
343 homologous TFs that exert the same functions usually use similar motifs (e.g Pou2f1/POU2F1,  
344 Otx1/OTX1 and/or Foxa3/FOXA3). Interestingly though, we found that this is not the case for  
345 all TF motifs. We found that the human RHOXF1 motif (Figure 2b) is likely not utilized by Rhox  
346 proteins in mice even though more than 30 Rhox genes exist. Evidently, throughout multiple  
347 duplications, Rhox genes seem to have obtained other functionalities in mouse <sup>38</sup> in  
348 comparison to the two human RHOX genes that are expressed in reproductive tissues <sup>39</sup>.  
349 Therefore, although we found the human RHOXF1 motif to be highly active in mice, this motif  
350 is most likely utilized by other proteins such as the mouse specific Obox proteins. In support  
351 of this conclusion, expression patterns of Obox proteins appear to be tightly regulated during  
352 PD <sup>40</sup> (<sup>21</sup>). High expression of Obox 1/2/5/7 is observed from the zygote to 4C stage, while  
353 Obox3/6/8 are expressed and peak at later stages (Supp. Table 4). Notably, there is a  
354 significant sequence similarity of the homeobox domains but not in the other parts of the  
355 RHOXF1 and Obox protein sequences, which supports the similarity in binding specificity.  
356 Although the potential functional overlap of RHOXF1 and Obox factors remains unresolved,  
357 our inter-species analysis suggests an unappreciated function of these factors and their  
358 targets during PD, warranting an in depth investigation.

359 In the context of TF target prediction, the power of TOBIAS was particularly highlighted by the  
360 fact that the analysis could identify almost all known Dux targets. In addition to coding genes,  
361 our analysis disclosed novel Dux binding sites and significant footprint scores at LINE1  
362 encoding genomic loci, which appear to be activated at the 4C/8C stage. This finding is  
363 especially interesting because a recent study has shown that LINE1 RNA can interact with  
364 Nucleolin and Kap1 to repress Dux expression <sup>41</sup>. Therefore, our findings give rise to a kinetics  
365 driven model in which Dux not only initiates ZGA but also regulates its own termination by a

366 temporally delayed negative feedback loop. Exactly how this feedback loop is controlled  
367 remains to be determined.

### 368 **Limitations and outlook of footprinting analysis**

369 Despite the striking capability of DGF analysis, some limitations and dependencies of this  
370 method still remain. Amongst these is the need of high-quality TF motifs for matching footprint  
371 scores to individual TFs with high confidence. In other words, while the binding of a TF might  
372 create an effect that can be interpreted as a footprint, without a known motif, this effect cannot  
373 be matched to the corresponding TF. This becomes evident in the context of DPPA2/4, a TF  
374 described by several groups to act in PD and even upstream of Dux<sup>34</sup>. DPPA2/4 targets GC  
375 rich sequences<sup>34</sup>, but its canonical binding motif remains unknown. It also needs to be noted  
376 that footprinting analysis cannot take effects into account that arise from heterogeneous  
377 mixtures of cells wherein TFs are bound in some cells and in others not. Therefore, if not  
378 separated, the classification of differential binding will be an observation averaged across  
379 many cells, possibly masking subpopulation effects. Recent advances have enabled the  
380 application of ATAC-seq in single cells<sup>42</sup>, but this generates sparse matrices, rendering  
381 footprinting approaches on single cells elusive. However, we speculate that by creating  
382 aggregated pseudo-bulk signals from large clustered SC ATAC datasets, DGF analysis might  
383 also become possible in single cells.



## 384 **Conclusions**

385 Here, we have illustrated the TOBIAS framework as a versatile tool for ATAC-seq footprinting  
386 analysis which helps to unravel transcription factor binding dynamics in complex experimental  
387 settings that are otherwise difficult to investigate. We showed that entire networks of TF  
388 binding, which have previously been explored using a combination of omics methods, can be  
389 recapitulated to a great extent by DGF analysis, which requires only ATAC-seq and TF motifs.  
390 From a global perspective, we provided new insights into PD by quantifying the stage-specific  
391 activity of specific TFs. Furthermore, we highlighted the usage of TOBIAS to study specific  
392 transcription factors as exemplified by our investigations on Dux. Finally, we used the specific  
393 TF target predictions to gain insights into the local binding dynamics of Dux in the context of  
394 TF-activation networks, repeat regions and novel genetic elements.

395 In conclusion, we present TOBIAS as the first comprehensive software that performs all steps  
396 of DGF analysis, natively supports multiple experimental conditions and performs visualization  
397 within one single framework. Although we utilized the process of PD as a proof of principle,  
398 the modularity and universal nature of the TOBIAS framework enables investigations of  
399 various biological conditions beyond PD. We believe that continued work in the field of DGF,  
400 including advances in both software and wet-lab methods, will validate this method as a  
401 resourceful tool to extend our understanding of a variety of epigenetic processes involving TF  
402 binding.

## 403 **Declarations**

### 404 **Ethics approval and consent to participate**

405 Not applicable.

### 406 **Consent for publication**

407 Not applicable.

### 408 **Availability of data and materials**

409 The TOBIAS software is available on GitHub at: <https://github.com/loosolab/TOBIAS>.

410 Excerpts of the data analyzed here are accessible for dynamic visualization at:

411 <http://loosolab.mpi-bn.mpg.de/tobias-meets-wilson>. All raw data analyzed are available from

412 GEO or ENCODE as described in Methods. The complete TOBIAS output for the analysis of

413 the Dux overexpression dataset can be downloaded from:

414 [https://figshare.com/projects/Digital\\_Genomic\\_Footprinting\\_Analysis\\_of\\_ATAC-](https://figshare.com/projects/Digital_Genomic_Footprinting_Analysis_of_ATAC-)

415 [seq\\_dataset\\_from\\_preimplantation\\_timepoints\\_via\\_TOBIAS/69959](https://figshare.com/projects/Digital_Genomic_Footprinting_Analysis_of_ATAC-seq_dataset_from_preimplantation_timepoints_via_TOBIAS/69959).

### 416 **Competing interests**

417 None to declare

### 418 **Funding**

419 This work was funded by the Max Planck Society, the German Research Foundation (DFG),

420 grant KFO309 (project number 284237345, epigenetics core unit) to ML, and by the Cardio-

421 Pulmonary Institute (CPI), EXC 2026, Project ID: 390649896 to ML.

### 422 **Authors' contributions**

423 MB, CK, JK and ML wrote the manuscript. MB, PG, HS, AP, KK, RW, AF and JP performed

424 the bioinformatics analysis. JK, TB and ML directed, coordinated and supervised the work.

## 425 **Acknowledgements**

426 We would like to thank the IT-group at MPI-BN for continued support with IT-infrastructure.

427 We would also like to thank Marius Dieckmann, the administrator of the Kubernetes cluster

428 in Gießen, for his support and help in implementing the TOBIAS-Nextflow Cloud version.

## 429 **Methods**

### 430 **Datasets**

<b>Organism</b>	<b>Deposited data</b>	<b>Source</b>	<b>Identifier</b>
Mouse	ATAC-seq, RNA-seq and ChIP-seq from mESC control and Dux overexpression	<sup>5</sup>	GEO: GSE85632
Mouse	ATAC-seq and RNA-seq from various preimplantation stages	<sup>21</sup>	GEO: GSE66390
Human	ATAC-seq and RNA-seq from various preimplantation stages	<sup>20</sup>	GEO: GSE101571

431

432 For all public data sets used in this study (see table above), raw files were obtained from the  
433 European Nucleotide Archive <sup>43</sup> and processed as described in the methods section. See also  
434 methods section “Comparison of TOBIAS to existing methods” for links to the ENCODE data  
435 used for method validation.

### 436 **Processing of ATAC-seq data**

437 Raw sequencing fastq files were assessed for quality, adapter content and duplication rates  
438 with FastQC v0.11.7, trimmed using cutadapt <sup>44</sup> and aligned with STAR v2.6.0c <sup>45</sup> (parameters:  
439 “--alignEndsType EndToEnd --outFilterMismatchNoverLmax 0.1 --  
440 outFilterScoreMinOverLread 0.66 --outFilterMatchNminOverLread 0.66 --outFilterMatchNmin  
441 20 --alignIntronMax 1 --alignSJDBoverhangMin 999 --alignEndsProtrude 10 ConcordantPair -

442 -alignMatesGapMax 2000 --outMultimapperOrder Random --outFilterMultimapNmax 999 --  
443 outSAMmultNmax 1”) to either the mouse or human genome using Mus\_musculus.GRCm38  
444 or Homo\_sapiens.GRCh38 versions from Ensembl <sup>46</sup>. Accessible regions were identified by  
445 peak calling for each sample separately using MACS2 (parameters: “--nomodel --shift -100 --  
446 extsize 200 --broad”) <sup>47</sup>. Peaks from each sample were merged to a set of union peaks across  
447 all conditions using “bedtools merge”. Each union peak was annotated to the transcriptional  
448 start site of genes (GENCODE <sup>48</sup>) in a distance of -10000/+1000 from the TSS using UROPA  
449 <sup>49</sup>.

## 450 Processing of RNA-seq data

451 Raw reads were assessed for quality, adapter content and duplication rates with FastQC  
452 v0.11.7, trimmed using cutadapt <sup>44</sup> and aligned with STAR v2.6.0c <sup>45</sup> (parameters: “--  
453 outFilterMismatchNoverLmax 0.1 --outFilterScoreMinOverLread 0.9 --  
454 outFilterMatchNminOverLread 0.9 --outFilterMatchNmin 20 --alignIntronMax 200000 --  
455 alignMatesGapMax 2000 --alignEndsProtrude 10 ConcordantPair --outMultimapperOrder  
456 Random --outFilterMultimapNmax 999”) to either the mouse or human genome using  
457 Mus\_musculus.GRCm38 or Homo\_sapiens.GRCh38 versions from Ensembl <sup>46</sup>. Differentially  
458 expressed genes were identified using DESeq2 v1.22 <sup>50</sup>. Only genes with a minimum log<sub>2</sub> fold  
459 change of  $\pm 1$ , a maximum Benjamini–Hochberg corrected P-value of 0.05 and a minimum  
460 combined mean of five reads were classified as significantly differentially expressed.

## 461 Processing of ChIP-seq data

462 Raw sequencing files in fastq format were quality assessed by Trimmomatic by trimming reads  
463 after a quality drop below a mean of Q15 in a window of 5 nucleotides <sup>51</sup>. All reads longer than  
464 15 nucleotides were aligned versus the mouse genome version mm10, keeping just unique  
465 alignments (parameters: --outFilterMismatchNoverLmax 0.2 --outFilterScoreMinOverLread  
466 0.66 --outFilterMatchNminOverLread 0.66 --outFilterMatchNmin 20 --alignIntronMax 1 --  
467 alignSJDBoverhangMin 999 --outFilterMultimapNmax 1 --alignEndsProtrude 10

468 ConcordantPair) by using the STAR mapper <sup>45</sup>. Read deduplication was done by Picard  
469 (<http://broadinstitute.github.io/picard/>).

## 470 Processing of transcription factor motifs

471 TF motifs were downloaded from JASPAR CORE 2018 <sup>52</sup>, the JASPAR PBM HOMEO  
472 collection and Hocomoco V11 <sup>53</sup> databases. We further included the human ARGFX\_3 motif  
473 from footprintDB <sup>54</sup> which originates from a HT-SELEX assay <sup>55</sup>. In annotation to the Dux/Dux4  
474 motifs of JASPAR and Hocomoco, we also included two TF motifs for MDUX/DUX4 created  
475 using MEME-ChIP <sup>56</sup> with standard parameters on the ChIP-seq peaks of <sup>35</sup> (GSE87279).

476 JASPAR motifs were linked to Ensembl gene ids by mapping the provided “Uniprot id” to the  
477 “Ensembl gene id” through biomaRt <sup>57</sup>. Hocomoco motifs were likewise linked to genes  
478 through the provided HGNC/MGI annotation. Due to the redundancy of motifs between  
479 JASPAR and Hocomoco, we further filtered the TF motifs to one motif per gene, preferentially  
480 choosing motifs originating from mouse/human respectively. For each TOBIAS run, we  
481 created sets of expressed TFs as estimated from RNA-seq in the respective conditions. This  
482 amounted to 590 motifs for the dataset on human preimplantation stages, 464 motifs for the  
483 dataset on mouse preimplantation, and 459 for the DuxOE dataset.

## 484 Maternal genes

485 Maternal genes for human and mouse were downloaded from the REGULATOR database <sup>22</sup>.  
486 Entrez gene ids were converted to Ensembl gene ids using biomaRt <sup>57</sup> and subsequently  
487 matched to available TF motifs as previously explained.

## 488 Overlap of Dux binding sites to repeat elements

489 Repeat elements for mm10 were downloaded from UCSC  
490 (<http://hgdownload.soe.ucsc.edu/goldenPath/mm10/database/rmsk.gz>). Overlap of Dux sites  
491 to individual repeat elements (as seen in figure 4G) was performed using “Bedtools intersect”.  
492 The sum of overlaps were counted by repeat class (LINE1/LTR).

## 493 Visualization

494 All TF-score heatmaps were generated by R Version 3.5.3 and complex heatmap package  
495 version 3.6<sup>58</sup>. Individual gene views were generated by loading TOBIAS output tracks into  
496 IGV version 2.6.2<sup>59</sup> or using the TOBIAS PlotTracks module, which is a wrapper for the  
497 svist4get visualization tool<sup>60</sup>. TF networks were drawn with Cytoscape version 3.7.1<sup>61</sup>.  
498 Heatmaps of genomic signal density were generated using Deeptools version 3.3.0<sup>62</sup>. All  
499 other figures, such as footprint plots, volcano plots and motif clustering dendrograms were  
500 generated by the TOBIAS visualization modules as described below.

## 501 The TOBIAS framework

502 In developing TOBIAS, we found that there were six main areas of DGF which had not been  
503 comprehensively addressed in the context of ATAC-seq footprinting analysis:

- 504 ● All-in-one framework including bias correction, footprinting, quantification of protein  
505 binding and visualization
- 506 ● Investigation of TF binding on a global level (which TFs are more bound globally ) as  
507 well as the locus-specific level (which TF binds to which genomic locations including  
508 statistics on differential binding)
- 509 ● Consideration of the redundancy and similarity of known TF binding motifs in the  
510 context of footprinting
- 511 ● A scoring model for TF-DNA binding taking into account the potential lack of a  
512 canonical footprint effect
- 513 ● Comparison and quantification of TF binding activity within complex experimental  
514 settings (multiple conditions or time series)
- 515 ● Automated workflows for recurring analysis tasks

516

517 Modules enabling these individual analysis steps are included in the TOBIAS package, which  
518 is publicly available at Github (<https://github.com/loosolab/TOBIAS>) as well as on PyPI and

519 Bioconda. Besides the examples given in the repository README, we also provide a Wiki  
520 (<https://github.com/loosolab/TOBIAS/wiki>) which introduces some of the individual software  
521 modules. We used the pre-defined workflows in Snakemake and NextFlow to run the full  
522 analysis. The single modules are explained in more detail below.

### 523 **Bias correction (TOBIAS ATACorrect module)**

524 Each Tn5-cut site is defined as the 5' end of the read shifted by +5 at the plus strand and -4  
525 at the minus strand to center the transposase event. Using the mapped reads from closed  
526 chromatin, ATACorrect builds a dinucleotide weight matrix<sup>63</sup> representing the preference of  
527 Tn5 insertion. In contrast to the classical position weight matrix (PWM) the dinucleotide weight  
528 matrix (DWM) captures the inter-base relationships which arise due to the palindromic nature  
529 of the bias. A background model is similarly built by shifting all reads +100bp as described by  
530<sup>64</sup>.

531 Reads within open chromatin peaks are then corrected by estimating the expected number of  
532 cuts per base pair and subtracting this from the observed cut sites as follows (modified from  
533<sup>65</sup>):

$$534 \quad c_i = x_i - e_i$$

535 where

$$536 \quad e_i = \hat{x}_i * \hat{b}_i, \quad \hat{x}_i = \sum_{j=i-50}^{i+50} x_j, \quad \hat{b}_i = \frac{b_j}{\sum_{j=i-50}^{i+50} b_j}$$

537 where  $x_i$  is the observed number of cuts,  $e_i$  is the expected number of cuts,  $b_i$  is the calculated  
538 bias level, and  $c_i$  is the corrected number of cuts at position  $i$ . To limit the influence of low-bias  
539 positions in the calculation of  $\hat{b}_i$ , a lower limit is set for  $b_i$  by calculating the fit of cutsites vs.  
540 bias to a rectified linear unit function (ReLU) in moving 100bp-windows and setting every  $b_i$   
541 below the linear fit to 0. This calculation is performed for all base pairs within open chromatin,

542 setting all other positions to 0. Lastly, each  $c_i$  is rescaled to fit the original sum of cuts  $\widehat{x}_i$   
543 for each window.

#### 544 **Footprinting (TOBIAS ScoreBigwig module)**

545 We estimate footprint scores across open chromatin regions by calculating:

$$546 \quad FP = \bar{x}_{flank} - \bar{x}_{mid}$$

547 where

$$548 \quad \bar{x}_{flank} = \frac{\sum_{i=j}^{j+wf} x_i + \sum_{i=j+wf+wm}^{j+2*wf+Wm} x_i}{2 * wf} \text{ for } x_i > 0$$

$$549 \quad \bar{x}_{mid} = \frac{\sum_{i=j+wf}^{j+wf+wm} x_i}{wm} \text{ for } x_i < 0$$

550  $x_i$  is the number of cuts at position  $i$ ,  $wf$  = width of flank in bp,  $wm$  = width of middle (footprint)  
551 in bp. The defaults used are:  $wf = [10;30]$ ,  $wm = [20;50]$ .

552 The term  $\bar{x}_{mid}$  will be negative and will therefore raise the score if there is a high depletion of  
553 cuts in the footprint (middle). If there is no depletion, the score will simplify to the mean of cuts  
554 in the flanking regions, representing accessibility. It is therefore not necessary to see a  
555 canonical footprint shape for the footprint score to be high. The footprint score can be  
556 interpreted as higher scores being more evidence that a protein was bound at a given position.

#### 557 **Estimation of transcription factor states and pairwise comparison between conditions** 558 **(TOBIAS BINDetect module)**

559 To match the calculated footprint scores to potential binding sites, TOBIAS BINDetect  
560 integrates genomic sequence, footprint scores from several conditions and motifs to identify  
561 up- and down regulated TFs based on footprint scores.

562 In the first step of the algorithm, the MOODS library (<https://github.com/jhkorhonen/MOODS>  
563 <sup>66</sup>) is used to detect TF binding sites (within peaks) with a p-value threshold of 1e-4.



564 Background base pair probabilities are estimated from the input peak set. Subsequently, each  
565 binding site is matched to footprint scores for each condition. Simultaneously, a background  
566 distribution of values is built by randomly subsetting peak regions at ~200bp intervals, and the  
567 scores from each condition are normalized to each other using quantile normalization. These  
568 values are used to calculate a distribution of background log2FCs for each pairwise  
569 comparison of conditions.

570 Overlaps between the TFBS identified in the first step are quantified by creating a distance  
571 matrix of TFs. The distance between a TF pair (TF1;TF2) is calculated as:

572

$$573 \text{dist}_{TF1;TF2} = 1 - \max(\text{overlap}_{TF1;TF2} / \text{total}_{TF1}, \text{overlap}_{TF2;TF1} / \text{total}_{TF2})$$

574

575 where  $\text{total}_{TF1}$  and  $\text{total}_{TF2}$  are the total base pairs of all TF1 and TF2 sites respectively  
576 and  $\text{overlap}_{TF1;TF2}$  is the amount of base pairs of TF1 which overlap with TF2 sites. The  
577 max-statement ensures that the overlap is calculated with regards to the shortest TF motif.

578 In the second step of the algorithm, every TF binding site found (for each motif given as input)  
579 is split into bound and unbound sites based on a score threshold per condition. The threshold  
580 is set at the level of significance of a normal-distribution fit to the background distribution of  
581 scores (user-defined p-value). As well as the per-condition split, each site is assigned a  
582 log2FC (fold change) per comparison, which represents whether the binding site has  
583 larger/smaller footprint scores in comparison. The global distribution of log2FC's per TF is  
584 compared to the background distributions to calculate a *differential binding score*, which is  
585 calculated as:

$$586 \frac{(\bar{x}_o - \bar{x}_b)}{((std_o + std_b) / 2)}$$

587 where  $\bar{x}_o$ ,  $std_o$  and  $\bar{x}_b$ ,  $std_b$  are the means and standard deviations of the observed and  
588 background log2FC distributions respectively. A p-value is also calculated by subsampling

589 100 log2FCs from the background and calculating the significance of the observed change  
590 (Python's `scipy.stats.ttest_1samp`). By comparing the observed log2FC distribution to the  
591 background log2FC, the effects of any global differences due to sequencing depth, noise etc.  
592 are controlled.

593 The differential binding scores and p-values are visualized as a volcano plot per condition-  
594 comparison. All TFs with  $-\log_{10}(\text{p-value})$  above the 95% quantile or differential binding scores  
595 smaller/larger than the 5% and 95% quantiles (top 5% in each direction) are colored and  
596 shown with labels. Below the plot, hierarchical clustering of the TFBS-distance matrix is shown  
597 and all TFs with distances less than 0.5 (overlap of 50% of bp) are colored as separate  
598 clusters.

599 The result of BINDetect is a folder-structure containing an overview of all potential binding  
600 sites (as .bed as well as excel-files), the predicted split into bound and unbound sites, and a  
601 global overview of differentially bound TFs per condition-comparison.

### 602 **Visualizing aggregate plots and calculation of footprint depth (TOBIAS PlotAggregate** 603 **module)**

604 Footprints are visualized using the subtool "TOBIAS PlotAggregate". Aggregate footprints are  
605 created by aligning genomic signals centered on all binding sites (taking into account  
606 strandedness), to create a matrix of ( $n$  sites)  $\times$  ( $n$  bp). The aggregate signal is calculated as  
607 the mean of each column (each bp). The default of  $\pm 60$ bp from the motif center was used  
608 throughout this manuscript.

609 The aggregate footprinting depth (FPD), which is applied in Supp. Figure 2c-d, was calculated  
610 for each TF as:

$$611 \quad FPD = \overline{signal}_{flank} - \overline{signal}_{middle}$$

612

613 where  $\overline{signal}_{middle}$  is the mean of the signal centered on the TFBS (30bp) and  
614  $\overline{signal}_{flank}$  is the mean of the signal in the remaining flanks ([-60;-15] and [+15;+60] bp)  
615 (See Supp. Figure 2b).  
616 Similarly to the investigations in previous literature <sup>16</sup>, we applied a mixture model from the  
617 Mixtools R package <sup>67</sup> to estimate the fractions of TFs with/without measurable footprints  
618 (Supp. Figure 2e).

619

## 620 **Transcription factor binding network (TOBIAS CreateNetwork module)**

621 The TF-TF network for Dux was built by subsetting all binding sites on the following  
622 characteristics: Bound in the promoter of a target gene, labeled “Unbound” in Control, labeled  
623 “Bound” in DuxOE, and log<sub>2</sub>FC footprint score increasing for DuxOE vs. Control. All targets  
624 were further reduced to only include genes encoding TFs with available motifs. Motifs were  
625 matched to genes as explained in the methods section “Processing of transcription factor  
626 motifs”. The network was then created using “TOBIAS CreateNetwork”. The result is a network  
627 of source and target nodes with directed edges, which in words can be described as: *Source*  
628 *TF* binds in the promoter of *Target TF*.

## 629 **TOBIAS framework output structure**

630 The output generated by the TOBIAS framework is organized in a hierarchical folder structure,  
631 which increases clarity of all steps of the analysis. The folder structure specifically organizes  
632 input data, pre-processing output like peak-calling and annotation, genomic tracks such as  
633 bias correction and footprints, as well as the local and global TF predictions. Particularly, the  
634 output for every individual TF investigated is arranged into separate folders containing TF  
635 specific plots, annotations and binding predictions. This structure makes it simple to use the  
636 output for further downstream analysis, as was showcased in this work. An exemplary output  
637 of the complete framework can be found at:

638 [https://figshare.com/projects/Digital\\_Genomic\\_Footprinting\\_Analysis\\_of\\_ATAC-](https://figshare.com/projects/Digital_Genomic_Footprinting_Analysis_of_ATAC-seq_dataset_from_preimplantation_timepoints_via_TOBIAS/69959)  
639 [seq\\_dataset\\_from\\_preimplantation\\_timepoints\\_via\\_TOBIAS/69959](https://figshare.com/projects/Digital_Genomic_Footprinting_Analysis_of_ATAC-seq_dataset_from_preimplantation_timepoints_via_TOBIAS/69959).

## 640 Validation

### 641 **Comparison of TOBIAS to existing methods**

642 Although footprinting tools for DNase-seq exist<sup>68-70 65, 71-73 74</sup>, not all can be applied to paired-  
643 end ATAC-seq data. We have focused our comparison on tools which are easily obtainable  
644 and installable, do not require ChIP-seq training-data, and are explicitly supporting ATAC-seq.  
645 We have additionally added two metrics for “Accessibility” and “PWM score” to compare  
646 TOBIAS to other footprinting-free metrics. The validation datasets and usage of existing tools  
647 are described in the following sections.

### 648 **Datasets**

649 The TOBIAS framework was benchmarked using ATAC-seq data from four human cell types:  
650 GM12878 (GEO: GSE47753), A549 (GEO: GSE114202), K562 (ENA: PRJNA288801) and  
651 HEPG2 (ENA: PRJEB30461). ATAC-seq data was trimmed using cutadapt<sup>44</sup> and mapped  
652 using Bowtie2<sup>75</sup>. All reads with a quality score <30 as well as non-proper paired reads were  
653 removed. All replicates were merged to one joined .bam-file of reads. Peaks were called using  
654 MACS2<sup>47</sup> with parameters "--nomodel --shift -100 --extsize 200 --broad --qvalue 0.01 --broad-  
655 cutoff 0.01". ChIP-seq peak regions (narrowPeak format) were downloaded from ENCODE  
656 and associated to motifs from Jaspar CORE 2018 using “MEME Centrimo”<sup>76</sup>. Only ChIP-seq  
657 experiments with motif enrichment > 1.0e-10 (Centrimo E-value) were kept. In case of more  
658 than one ChIP-seq experiment for the same target in the same cell type, the one with the  
659 highest motif enrichment was chosen. After filtering, there were 12 TFs for A549, 54 TFs for  
660 GM12878, 64 TFs for HepG2, and 87 TFs for K562 for a total of 217 ChIP-seq experiments  
661 matched to ATAC-seq. Bound binding sites per TF were defined as any TFBS within +/- 50bp  
662 from the paired ChIP-seq peak summit. In case of two or more binding sites per peak, the one  
663 closest to the summit was set to bound, and others were excluded from the analysis. Unbound

664 binding sites were defined as any TFBS not overlapping any ChIP-seq peak, as well as not  
665 overlapping bound sites from any other factors for this cell type. Bound and unbound sites  
666 were further filtered to only include TFBS falling within ATAC-seq peaks for the cell type in  
667 question.

## 668 **Bias correction approaches**

669 TOBIAS was compared to the existing bias correction methods as follows:

- 670 • **seqOutBias** <sup>(77)</sup>

671 The seqOutBias software was downloaded from GitHub  
672 (<https://github.com/guertinlab/seqOutBias>). Following the vignette for ATAC-seq,  
673 mappability files were created and ATAC-seq reads were corrected for plus/minus  
674 strand reads separately. After correction, we further shifted the positive and negative  
675 tracks +5 and -4bp respectively, as this was not performed by the tool itself.

- 676 • **HINT-ATAC** <sup>(14)</sup>

677 The HINT software was downloaded from PyPI as part of the RGT software suite. Bias-  
678 correction was performed from the ATAC-seq reads using the command “rgt-hint  
679 tracks --bc --bigWig <bam>”.

680

681 Aggregate footprints for each method across all (within peaks), bound and unbound binding  
682 sites (see explanation above) were visualized using “TOBIAS PlotAggregate”.

## 683 **Footprinting**

684 Existing footprinting tools were applied as follows:

- 685 • **msCentipede** <sup>(78)</sup>

686 The msCentipede software was downloaded from GitHub  
687 (<https://github.com/rajanil/msCentipede>). For each TF, the binding model was built  
688 using the 5000 TFBS with the highest PWM score genomewide. For model learning,  
689 the “--mintol” parameter was set to 1e-3 as a tradeoff between accuracy and speed.

690 The resulting models were then used to infer the posterior binding-probability of TFBS  
691 in peaks.

692 • **Wellington (70)**

693 The pyDNase software was downloaded from PyPI. Footprints in ATAC-seq peaks  
694 were estimated using “wellington\_footprints.py” with the “-A” option for ATAC-seq  
695 mode.

696 • **PIQ (79)**

697 The PIQ software was downloaded from Bitbucket ([https://bitbucket.org/thashim/piq-](https://bitbucket.org/thashim/piq-single/)  
698 [single/](https://bitbucket.org/thashim/piq-single/)). The script *bam2rdata.r* was used to bring the input .bam-file into the correct  
699 data format. Likewise, the script *pwmmatch.exact.r* was used to predict genomewide  
700 TFBS. Finally, footprinting scores for each TF were obtained using the script *perff.r* for  
701 each motif/cell type pair. The purity score was taken as the probability for a certain  
702 TFBS to be bound.

703 • **HINT-ATAC (14)**

704 The HINT software was downloaded from PyPI as part of the RGT software suite.  
705 Footprints were identified using the command “rgt-hint footprinting --atac-seq --paired-  
706 end --organism=hg38 <bam> <peaks>”. The output of HINT-ATAC footprinting is a  
707 .bed-file of footprint ranges ranked by tag count. All TFBS overlapping a footprint with  
708 more than 2/3 of the TFBS bases was assumed to be bound and scored using the tag  
709 count of the footprint. The rest of the TFBS (within peaks) were set to score 0 (low  
710 chance of protein binding). The auROC was calculated based on the ability of these  
711 scores to predict true protein binding. It should be noted that this affects the shape of  
712 the ROC curve, as all TFBS without overlaps are assumed to have the same probability  
713 of being bound. However, this is a characteristic of the method, and HINT-ATAC was  
714 therefore evaluated on the same premise as other tools.

715 • **Accessibility**

716 The “Accessibility” metric is defined as the sum of Tn5 insertions in a 300 basepair

717 window centered at the binding site. This score represents the accessibility of a certain  
718 region not taking into account local footprint information.

719 • **PWM score**

720 The score of the motif-sequence match at the specific TFBS. As this is based on  
721 sequence alone, the PWM-score is independent of chromatin accessibility.

722 Due to high computational times for some tools, the validation was limited to binding sites on  
723 human chromosome 1. On the basis of the ChIP-seq labels, the area under the ROC curve  
724 (auROC) was used to evaluate the predictive power of each method.

## 725 **Supplemental Information**

### 726 List of Supplementary Files

727 *Supplementary File 1: Visualization of different methods for Tn5 bias correction across 36 TFs*  
728 *with matched ChIP-seq. Each page contains footprints for a specific TF across all binding sites*  
729 *(in peaks), bound sites (overlapping ChIP-seq) and unbound sites (not overlapping ChIP-seq)*  
730 *for uncorrected/expected/corrected signals from different bias correction methods.*

731 *Supplementary File 2: The direct output file of the “TOBIAS BINDetect”-module containing*  
732 *differential binding plots across all pairwise-comparisons of human developmental stages.*

733 *Supplementary File 3: The direct output file of the “TOBIAS BINDetect”-module containing*  
734 *differential binding plots across all pairwise-comparisons of mouse developmental stages.*

735 *Supplementary File 4: The direct output file of the “TOBIAS BINDetect”-module containing*  
736 *differential binding plots between control (mESC) and DuxOE samples.*

737

## 738 List of Supplementary Tables

739 *Supplementary Table 1: Prediction of transcription factor binding across human*  
740 *2C/4C/8C/ICM/hESC clustered into co-active TFs. Each transcription factor is further linked to*  
741 *expression of the factor based on RNA-seq.*

742 *Supplementary Table 2: TOBIAS TF scores for human PD timepoints, correlated to*  
743 *corresponding RNA expression.*

744 *Supplementary Table 3: Prediction of transcription factor binding across mouse*  
745 *2C/4C/8C/ICM/mESC clustered into co-active TFs. Each transcription factor is further linked*  
746 *to expression of the factor based on RNA-seq.*

747 *Supplementary Table 4: Human and Mouse RNA expression for Obox and RHOX/Rhox genes*  
748 *during preimplantation developmental stages.*

749 *Supplementary Table 5: Full list of the predicted Dux binding sites as well as their change*  
750 *between mESC and DuxOE as predicted by TOBIAS.*



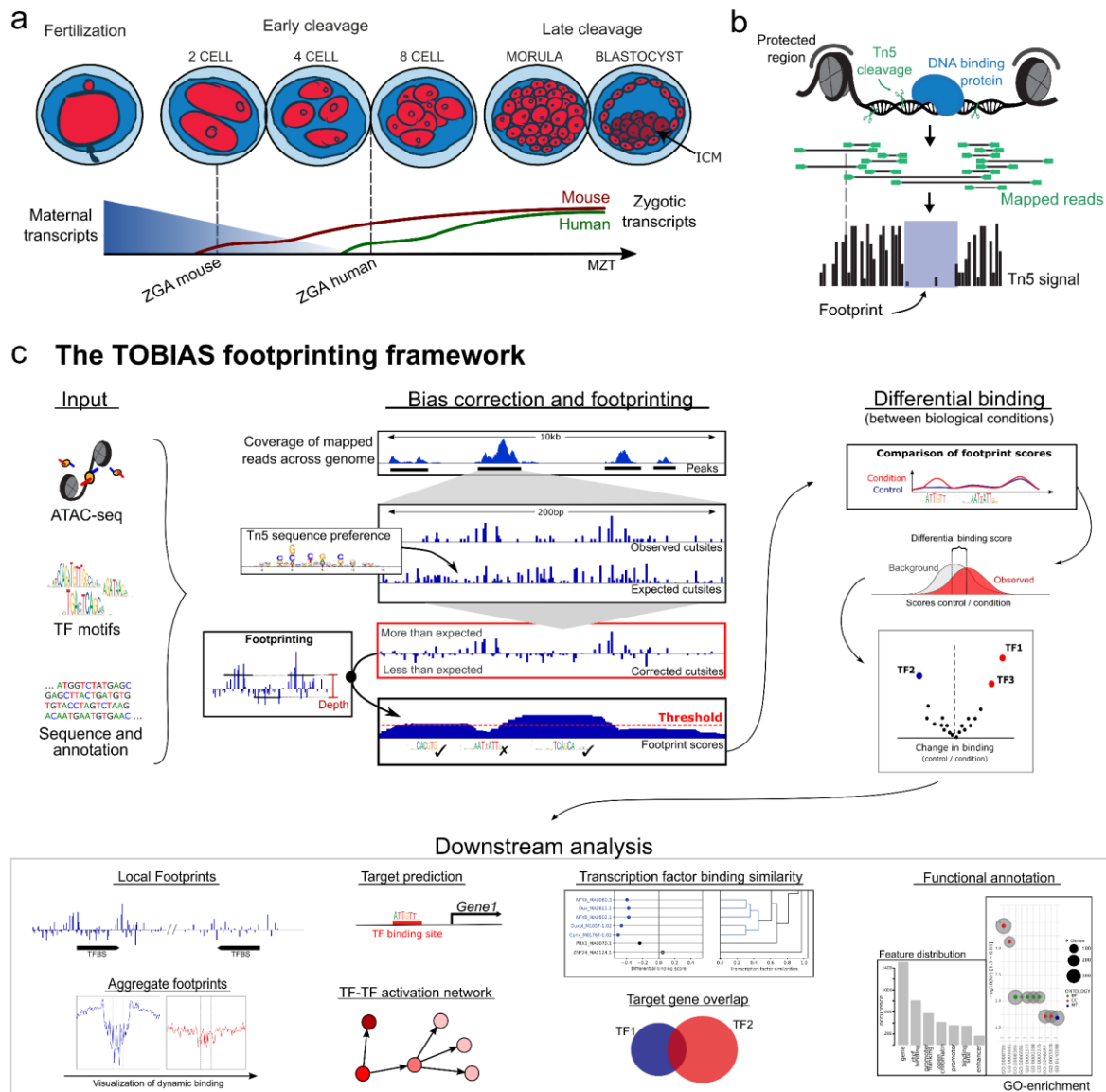
751 **Tables**

752 *Table 1: Comparison of features for ATAC-seq footprinting tools*

	<b>Footprinting tools for ATAC-seq</b>				
	<b>TOBIAS</b>	<b>HINT-ATAC</b>	<b>MsCentipede</b>	<b>PIQ</b>	<b>Wellington</b>
<b>Overview:</b>					
Year of publication	-	2019	2015	2014	2013
Tool availability	Github	Github	Github	Bitbucket	Github
Programming language	Python	Python	Python	R	Python
Type of footprinting (D: De novo, M: Motif-centric)	D	D	M	M	D
<b>Features:</b>					
Footprinting	✓	✓	✓	✓	✓
Tn5 bias-correction	✓	✓	✗	✗	✗
Size-adjustable footprinting algorithm	✓	✗	✗	✗	✗
Differential footprinting	✓	✓	✗	✗	✓
Time series footprinting (comparison of 2+ conditions)	✓	✗	✗	✗	✗
Calculation of TFBS (from motifs)	✓	✓	✗	✓	✗
TFBS clustering	✓	✗	✗	✗	✗
Consensus motifs for clustered motifs	✓	✗	✗	✗	✗
Output of genomic tracks	✓	✓	✗	✗	✓
Adjustable plotting of aggregate footprints	✓	✗	✗	✗	✗
Visualization of locus footprints	✓	✗	✗	✗	✗
Inference of TF-binding networks	✓	✗	✗	✗	✗
Predict bound/unbound state per TFBS	✓	✗	✓	✓	✗
<b>Usability:</b>					
Uses standard file formats	✓	✓	✗	✗	✓
Parallel computing	✓	✗	✗	✗	✓
Complete workflow available	✓	✗	✗	✗	✗
- Snakemake	✓	✗	✗	✗	✗
- Nextflow	✓				
Cloud computing supported	✓	✗	✗	✗	✗
Time to execute* (min)	7.2	46	2808	329	8
Package manager/installer	PyPI Bioconda	PyPI Bioconda	-	-	PyPI Bioconda

753 \* CPU time (using 30 cores if applicable) to perform bias correction (if applicable) and footprinting for  
754 GM12878 chromosome 1 using 54 transcription factors matched to ENCODE ChIP-seq.

## 755 Figures and figure legends



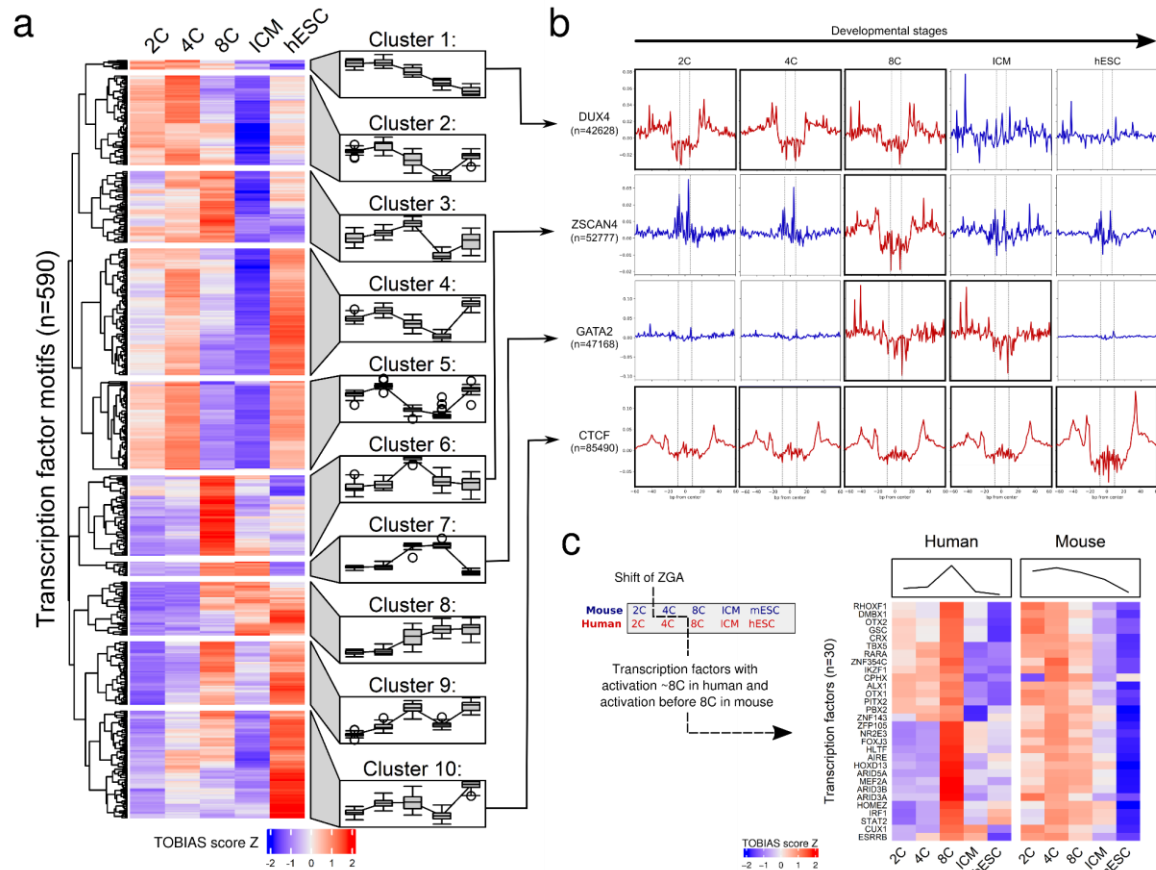
756

757 **Figure 1: The use of chromatin accessibility assays to investigate early developmental processes**

758 **(a) Early embryonic development in human and mouse.** The fertilized egg undergoes a series of divisions  
 759 ultimately creating the structure of the blastocyst. While maternal transcripts are depleted, the zygotic genome  
 760 is activated in waves. ZGA initiates in mouse at 2-cell stage and in human at the 4-8-cell stage.

761 **(b) The concept of footprinting using ATAC-seq.** The Tn5 transposase cleaves and inserts sequencing adapters  
 762 in open chromatin, but is unable to cut in chromatin occupied by e.g. nucleosomes or transcription factors. The  
 763 mapped sequencing reads are used to create a signal of single Tn5 insertion events, in which binding of  
 764 transcription factors is visible as depletion of signal (the footprint).

765 **(c) The TOBIAS digital genomic footprinting framework.** Using an input of sequencing reads from ATAC-seq,  
766 transcription factor motifs and sequence information, the TOBIAS footprinting framework detects local and  
767 global changes in transcription factor binding. Bias-correction of the Tn5 sequence preference enables detection  
768 of local chromatin footprints and matching to individual TFBS. Footprint scores are compared between conditions  
769 to define differentially bound TFs. The global binding map allows for a variety of downstream analysis such as  
770 visualization of local and aggregated footprints across conditions, prediction of target genes for each TF as well  
771 as comparison of binding specificity between several transcription factors. Functional annotation such as GO  
772 enrichment can be used to infer biological meaning of target gene sets.



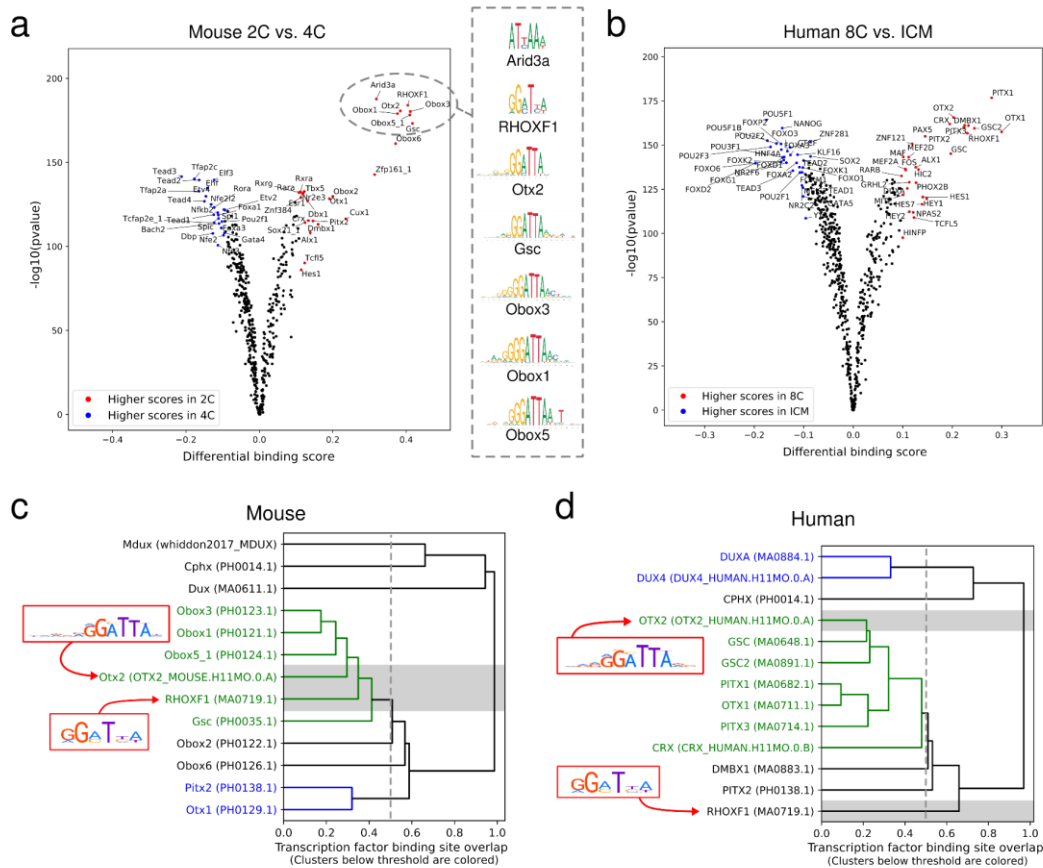
773

774 **Figure 2: TOBIAS enables investigation of global changes in transcription factor binding**

775 **(a) Clustering of transcription factor activities throughout development.** Each row represents one TF, each  
776 column a developmental stage; blue color indicates low activity, red color indicates high activity. In order to  
777 visualize cluster trends, each cluster is associated with a mean trend line and time point specific boxplots.

778 **(b) Bias-corrected ATAC-seq footprints reveal dynamic TF binding.** Aggregated footprinting plot matrix for  
779 transcription factor binding sites. Plots are centered around binding motifs (n=\* relates to the number of binding  
780 sites). Rows indicate TFs DUX4, ZSCAN4, GATA2, and CTCF; columns illustrate developmental stages from left to  
781 right. Active binding of the individual TFs is visible as depletion in the signal around the binding site (highlighted  
782 in red). Upper three TFs are related to developmental stages, CTCF acts as a universal control, generating a  
783 footprint in all conditions. See Supplementary Figure 3f for uncorrected footprints.

784 **(c) TF activity is shifted by ZGA onset in human and mouse.** Heatmaps show activity of known ZGA-related TFs  
785 for human (left) and mouse (right) across matched timepoints 2C / 8C / ICM / hESC (mESC). Mean TF activity (top  
786 panel) peaks at 4-8C stage in human and is shifted to 2-4C stage in mouse by the earlier ZGA onset.

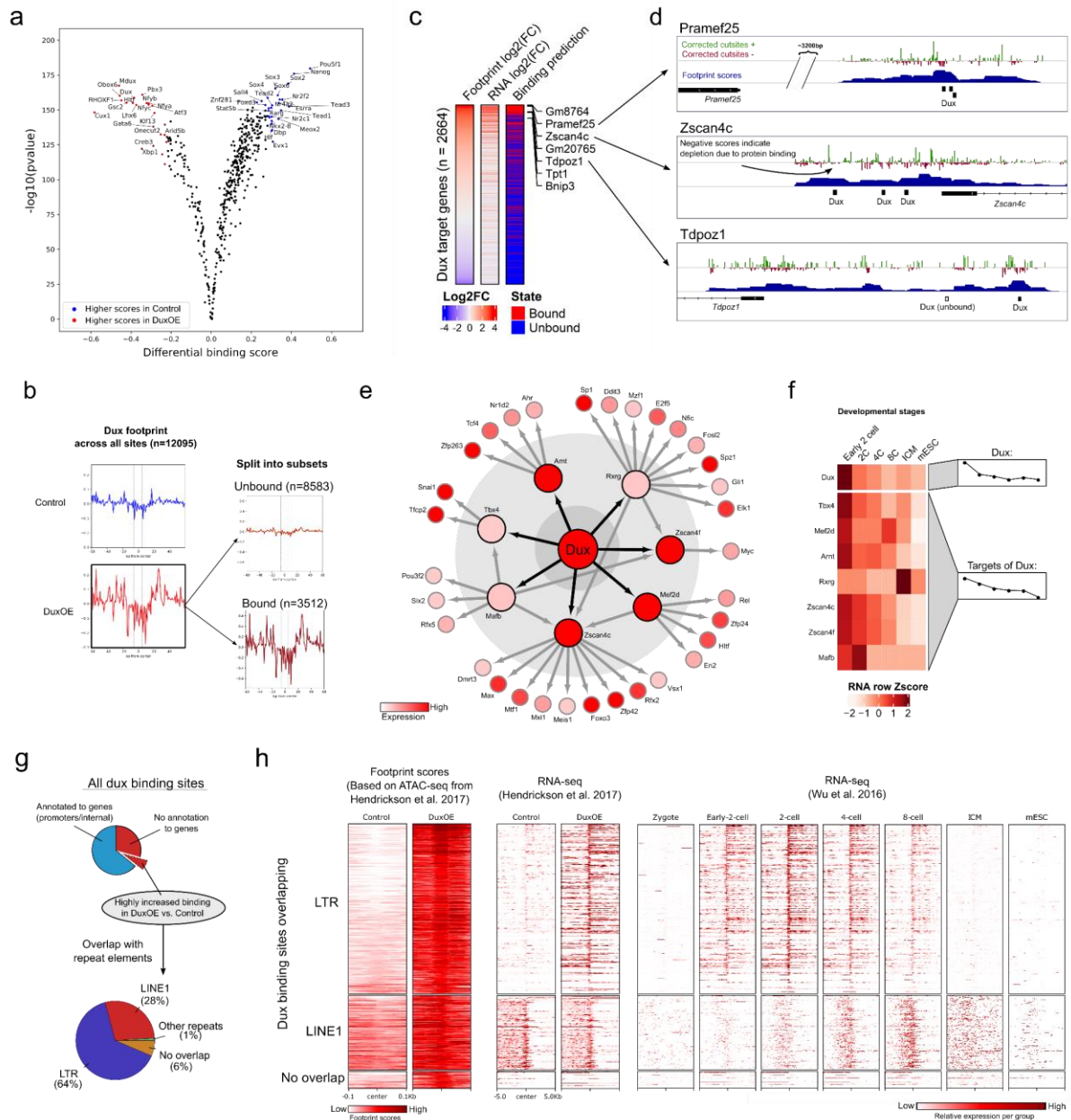


787

788 **Figure 3: Comparison of binding site overlaps shows specification of ZGA functions between mouse**  
 789 **and human**

790 *(a-b) Pairwise comparison of TF activity between developmental stages. The volcano plots show the differential*  
 791 *binding activity against the  $-\log_{10}(pvalue)$  (as provided by TOBIAS) of the investigated TF motifs; each dot*  
 792 *represents one motif. For (A) 2C stage specific/significant TFs are labeled in red, 4C specific factors are given in*  
 793 *blue. For (B) 8C stage specific/significant TFs are labeled in red, ICM specific factors are given in blue.*

794 *(c-d) Clustering of TF motifs based on binding site overlap. Excerpt of the global TF clustering based on TF*  
 795 *binding location, illustrating individual TFs as rows. The trees indicate genomic positional overlap of individual*  
 796 *TFBS with a tree-depth of 0.2 representing an overlap of 80% of motifs. Each TF is indicated by name and unique*  
 797 *ID in brackets. Clusters of TFs with more than 50% overlap (below 0.5 tree distance) are colored. (C) shows overlap*  
 798 *of motifs included in the mouse analysis, and (D) shows clustering of human motifs. Complete TF trees are*  
 799 *provided in Supp. Files 2 and 3.*



800

801 **Figure 4: Dux binding induces transcription at gene promoters and LTR sequences in mouse**

802 **(a) Volcano plot comparing TF activities between mDux GFP- (Control) and mDux GFP+ (DuxOE).** Volcano plot  
 803 showing the TOBIAS differential binding score on the x-axis and  $-\log_{10}(p\text{value})$  on the y-axis; each dot represents  
 804 one TF. DuxOE specific TFs are labeled in red and Control specific TFs are labeled in blue.

805 **(b) Aggregated footprint plots for Dux.** The aggregate plots are centered on the predicted binding sites for Dux  
 806 between Control and DuxOE condition. The total possible binding sites for DuxOE (n=12095) are separated into  
 807 bound and unbound sites (right). The dashed line represents the edges of the Dux motif.

808 **(c) Change in expression of genes near Dux binding sites.** The heatmap shows 2664 Dux binding sites found in  
 809 gene promoters. Footprint  $\log_2(FC)$  and RNA  $\log_2(FC)$  represent the changes between Control and DuxOE for

810 *footprints and gene expression, respectively.  $\log_2(FC)$  is calculated as  $\log_2(\text{DuxOE}/\text{Control})$ . The column “Binding*  
811 *prediction” depicts whether the binding site was predicted by TOBIAS to be bound/unbound in the DuxOE*  
812 *condition.*

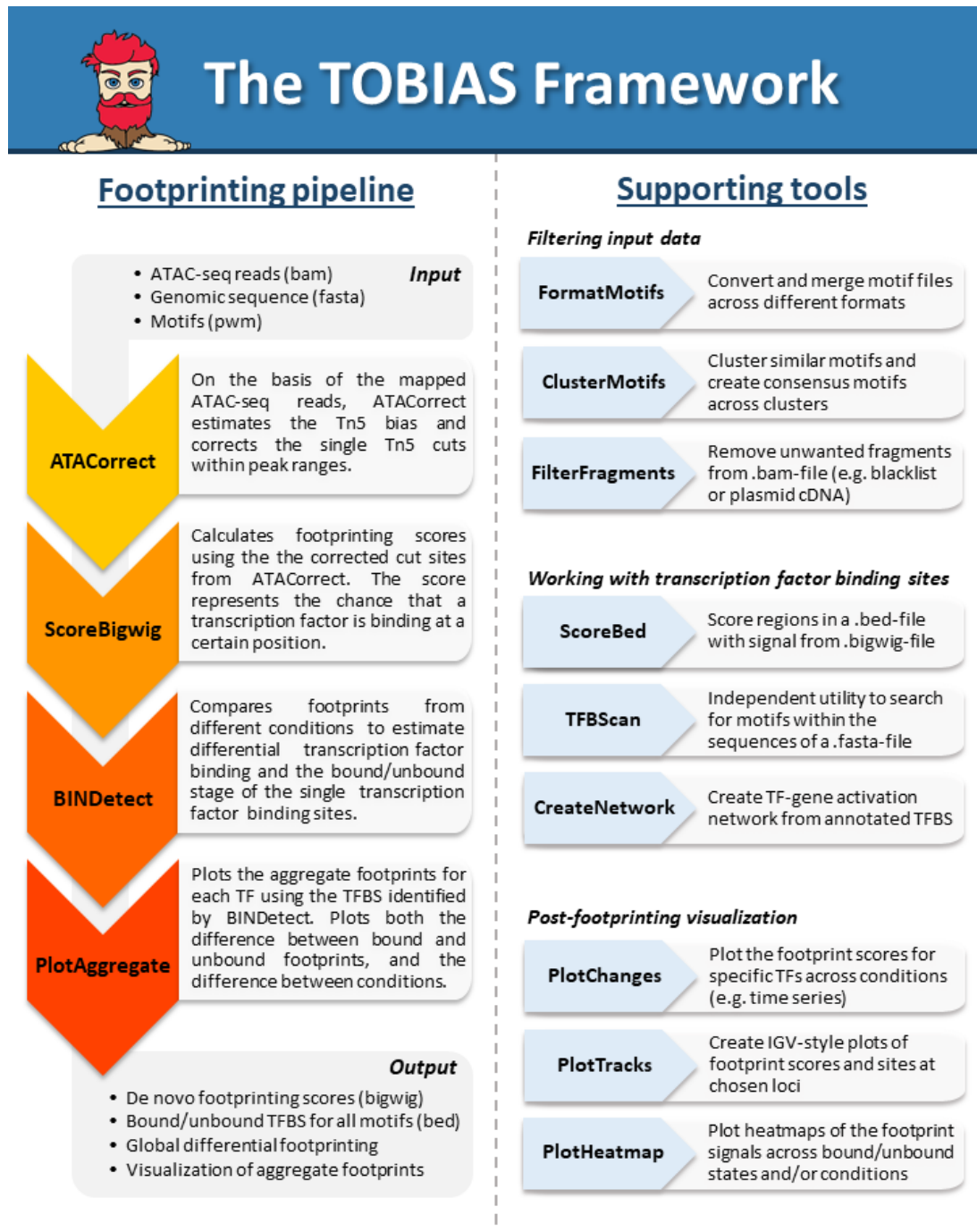
813 **(d) Genomic tracks showing footprint scores of Dux-binding.** *Genomic tracks indicating three DUX target gene*  
814 *promoters (one per row) and respective tracks for cut site signals (red/blue), TOBIAS footprints (blue), detected*  
815 *motifs (black boxes), and gene locations (solid black boxes with arrows indicating gene strand).*

816 **(e) Dux transcription factor network.** *The TF-TF network is built of all TFBS with binding in TF promoters with*  
817 *increasing strength in DuxOE ( $\log_2(FC)>0$ ). Sizes of nodes represent the level of the network starting with Dux*  
818 *(Large: Dux, Medium: 1st level, Small: 2nd level). Nodes are colored based on RNA level in the OE condition.*

819 **(f) Correlation of the Dux transcription factor network to expression during development.** *The heatmap depicts*  
820 *the in vivo gene expression during developmental stages from <sup>21</sup>. The right-hand group annotation highlights the*  
821 *difference in mean expression for each timepoint. The heatmap is split into Dux and target genes of Dux.*

822 **(g) Dux binding sites overlap with repeat elements.** *All potential Dux binding sites are split into sites either*  
823 *overlapping promoters/genes or without annotation to any known genes. The bottom pie chart shows a subset*  
824 *of the latter, additionally having highly increased binding ( $\log_2(FC)>1$ ), and overlapping LTR/LINE1 elements.*

825 **(h) Dux induces expression of transcripts specific for preimplantation.** *Genomic signals for the Dux binding sites*  
826 *which are bound in DuxOE with  $\log_2(FC)$  footprint score  $>1$  (i.e. upregulated in DuxOE) are split into overlapping*  
827 *either LTR, LINE1 or no known genetic elements (top to bottom). Footprint scores ( $\pm 100\text{bp}$  from Dux binding*  
828 *sites) indicate the differential Dux binding between control and DuxOE. RNA-seq shows the normalized read-*  
829 *counts from <sup>5</sup> and <sup>21</sup> within  $\pm 5\text{kb}$  of the respective Dux binding sites where red color indicates high expression.*

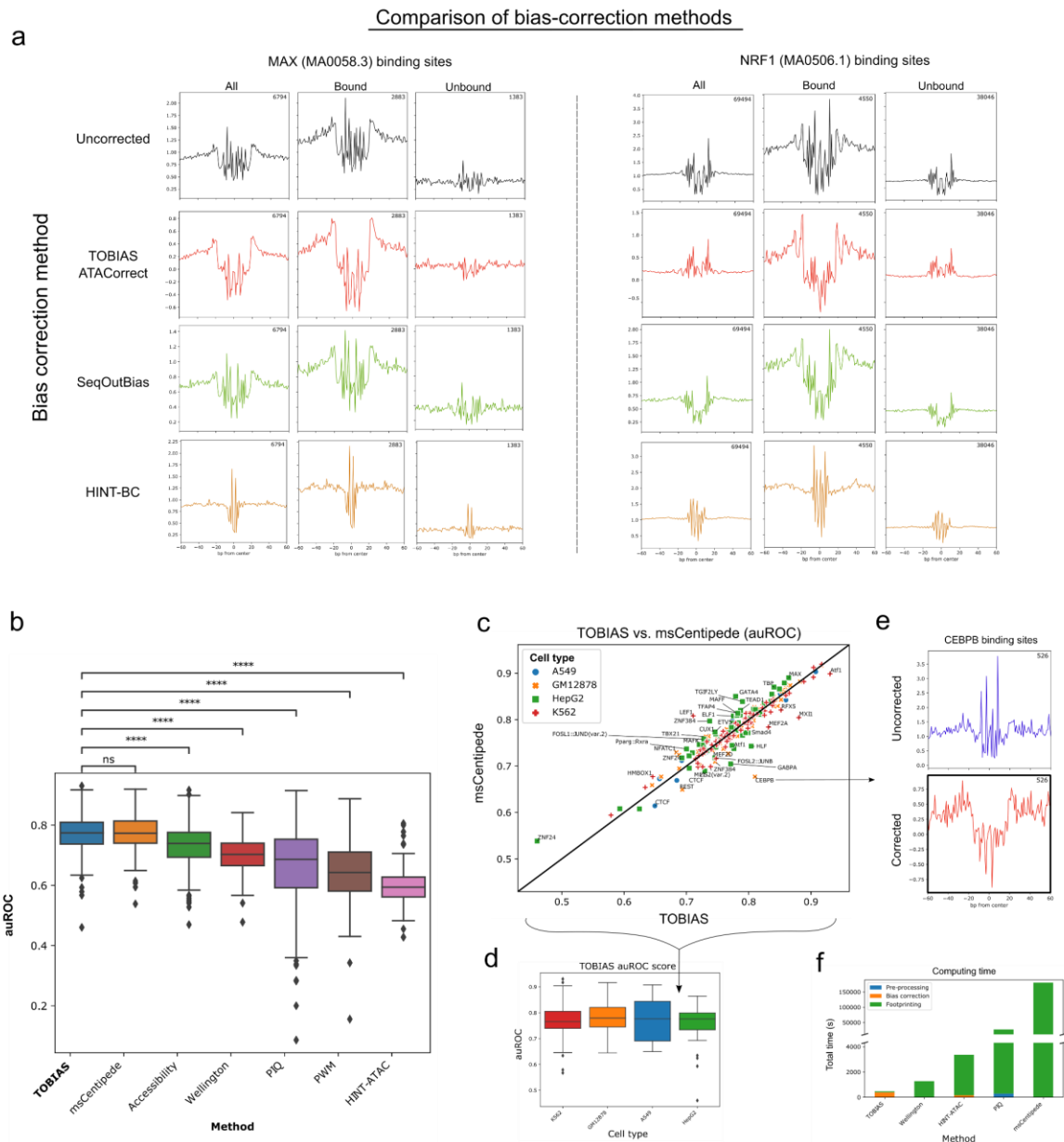


830

### 831 **Supplementary Figure 1: Overview of the TOBIAS framework tools**

832 *The TOBIAS tools are intended for use in a standardized pipeline as shown on the left. ATACorrect and*  
833 *ScoreBigWig corrects Tn5 cuts and calculates footprint scores respectively. Next, BINDetect introduces*  
834 *information from different transcription factor binding motifs to predict binding sites both within and across*  
835 *conditions. PlotAggregate can be used to visualize the single footprints. Furthermore, a large variety of*  
836 *supporting tools can be used at different stages of the pipeline, such as pre-filtering of .bam-files using*  
837 *FilterFragments or plotting of locus-specific footprints using PlotTracks.*





838

839 **Supplementary Figure 2: Comparison of existing bias-correction and footprinting methods**

840 **(a) Comparison of aggregate footprints for different bias-correction methods.** Bound and unbound

841 transcription factor binding sites for MAX and NRF1 are shown across uncorrected signal (pileup of Tn5

842 insertions), TOBIAS ATACCorrect, SeqOutBias and HINT-BC correction methods. An overview of all included TFs

843 from cell type GM12878 can be found in Supplementary File 1.

844 **(b) Comparison of predictive ability across different footprinting methods.** The auROC is calculated based on

845 ENCODE ChIP-seq for 217 TFs and compared across methods. Significance (Mann-Whitney U test, \*\*\* equals

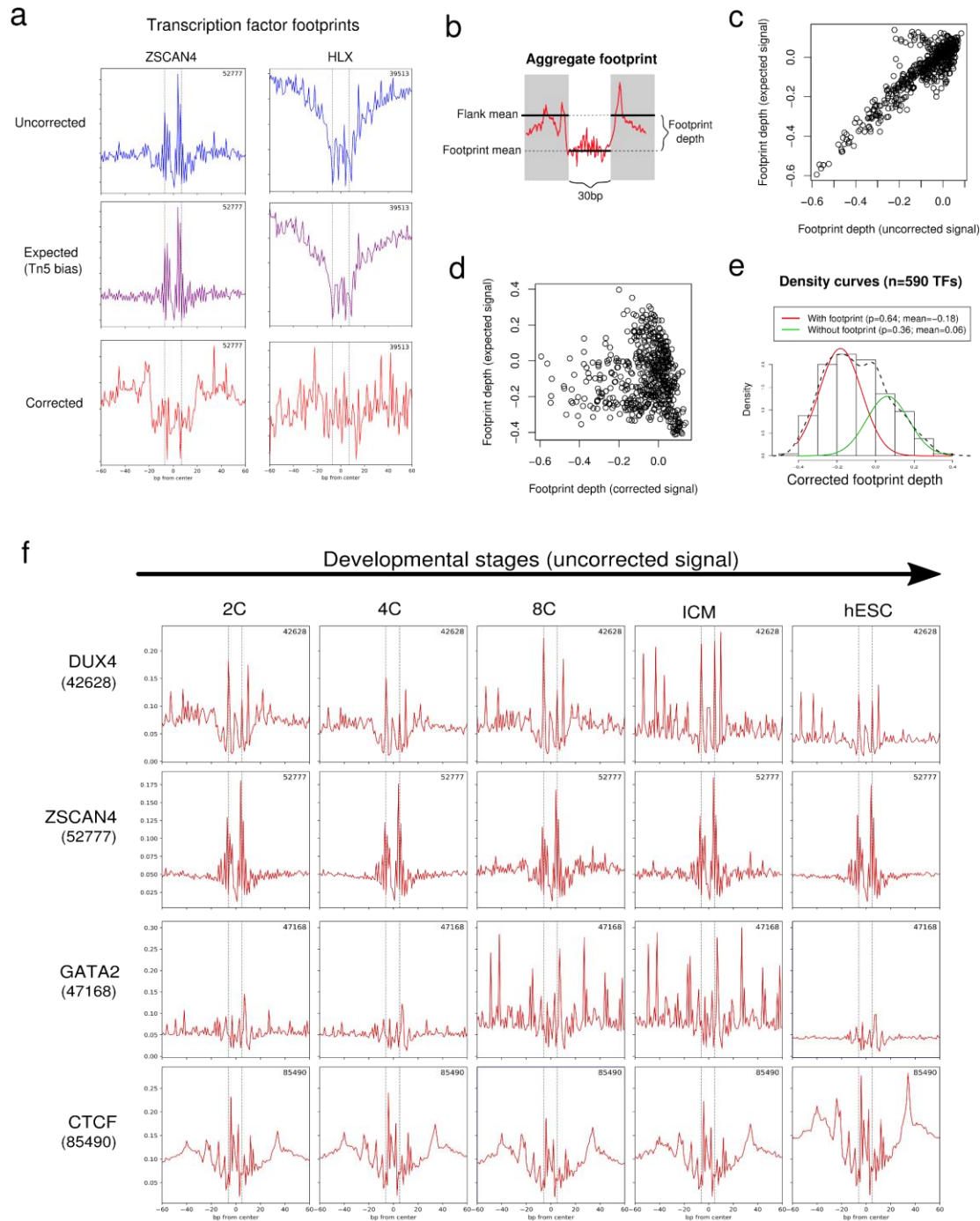
846  $p \leq 1.0e^{-4}$ ) is indicated as asterisk.

847 **(c) Scatterplot comparing the auROC of TOBIAS and msCentipede.** Each point represents one TF, which is  
848 colored and marked dependent on cell type. The diagonal line represents equal auROC between TOBIAS and  
849 msCentipede.

850 **(d) Validation of TOBIAS across cell types.** The auROC of TOBIAS predictions across cell types K562 (n=67),  
851 GM12878 (n=54), HepG2 (n=64) and A549 (n=11).

852 **(e) Aggregate footprints for CEBPB.** The aggregate footprints for true CEBPB binding sites (bound sites verified  
853 by ChIP-seq). Whereas the uncorrected ATAC-seq is insufficient to uncover a footprint, the corrected ATAC-seq  
854 signal exhibits a clear footprint for CEBPB binding sites.

855 **(f) Comparison of computing times for footprinting tools.** The CPU run time for each tool is measured across  
856 the three tasks of “pre-processing” (only for PIQ), “bias-correction” (only for TOBIAS and HINT-ATAC) and  
857 footprinting (all tools).



858

859 **Supplementary Figure 3: Tn5-bias correction is important for visualization of footprints from ATAC-**  
 860 **seq**

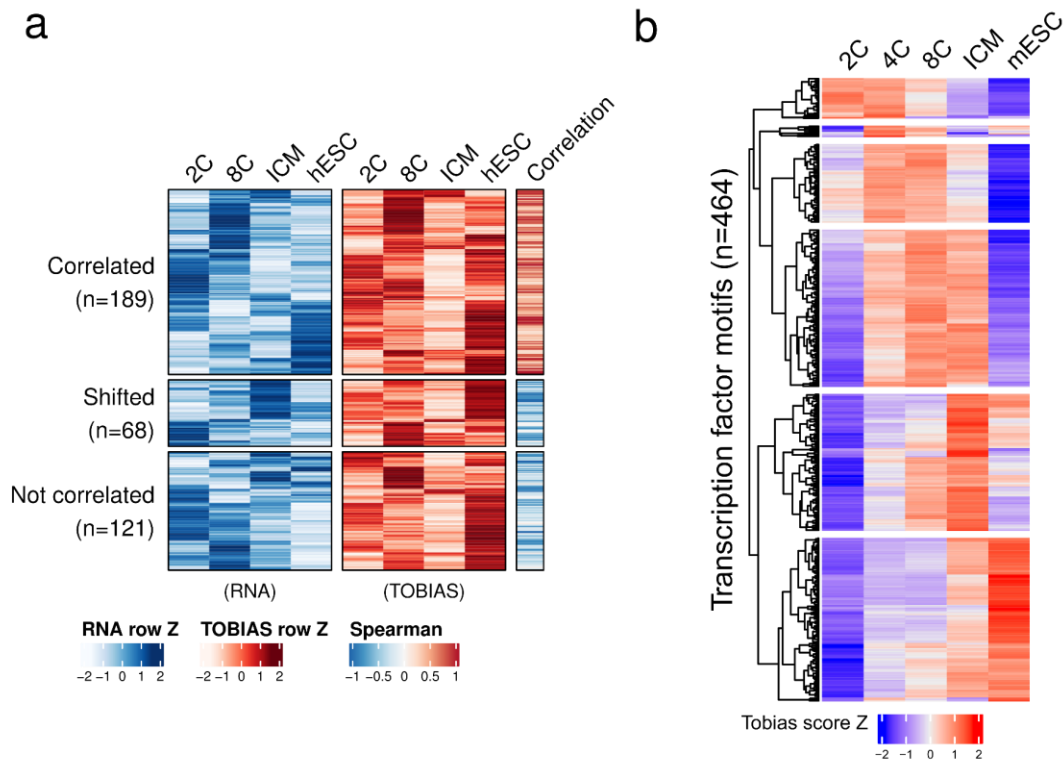
861 *(a) Examples of Tn5-bias correction using “expected”-intermediates. The figure shows the aggregate footprints*  
 862 *for transcription factors ZSCAN4 and HLX across the uncorrected, expected and corrected Tn5 signals. The*  
 863 *number in the right-hand corner represents number of binding sites included in the plot.*

864 *(b) Aggregate footprint depth model. The footprint depth is calculated using a similar metric as described in <sup>16</sup>.*

865 **(c-d) Uncorrected and corrected Tn5-bias.** The scatter plots show the correlation between depth of footprints  
866 for uncorrected vs. expected footprints (c) and corrected vs. expected footprints (d).

867 **(e) Mixture model of all footprinting depths.** The mixture model shows that 65% of motifs fall into the category  
868 of a measurable footprint in the aggregated profile. Data is based on 590 motifs in hESC.

869 **(f) A depiction of uncorrected footprint aggregates across time points for transcription factors DUX4, ZSCAN4,**  
870 **GATA2 and CTCF.** In contrast to the corresponding corrected signals seen in Figure 2A, the footprints are hardly  
871 visible in the uncorrected aggregates.

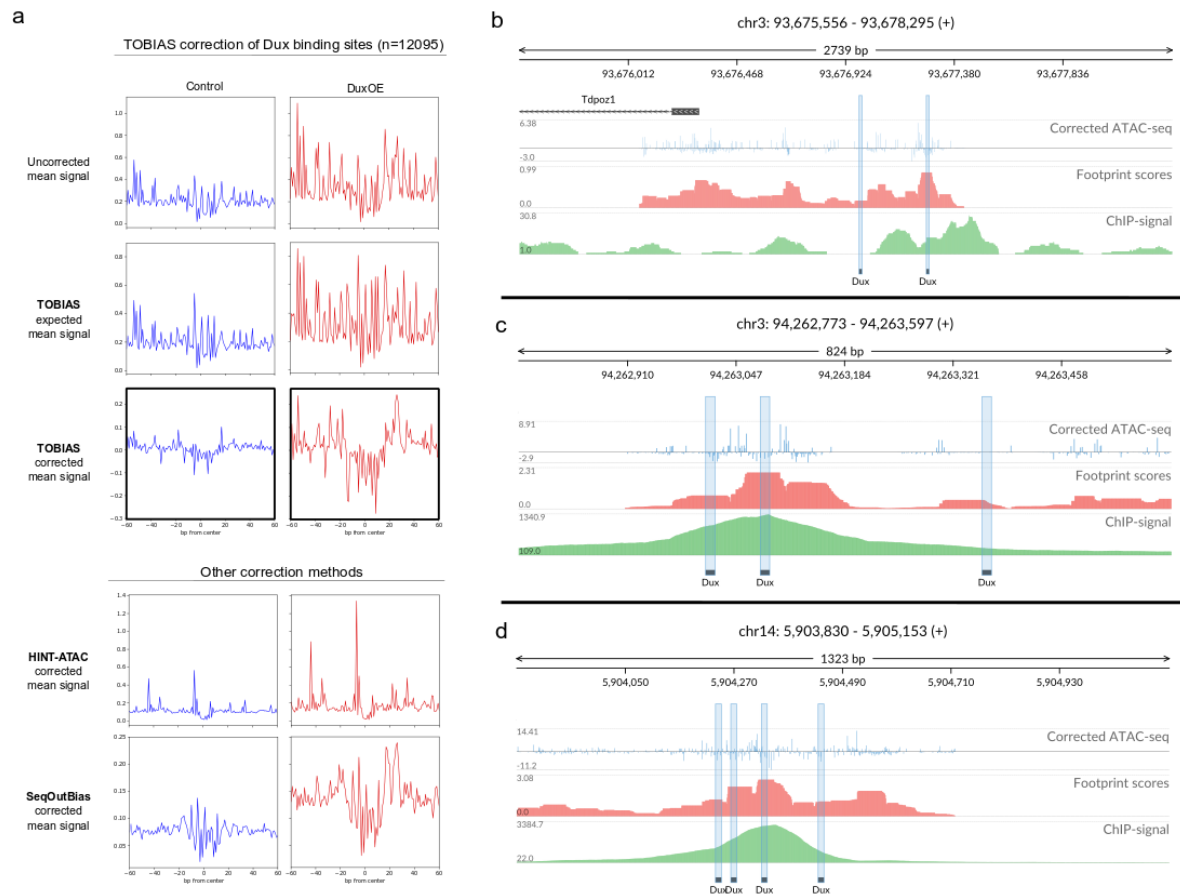


872

873 **Supplementary Figure 4: Transcription factor activity and expression during mouse and human**  
 874 **development**

875 **(a) Correlation of footprints and RNA-seq.** The left heatmap (blue) depicts expression of transcription factor  
 876 clusters in the respective human developmental stages. The left heatmap (red) depicts the corresponding TOBIAS  
 877 scores. Spearman column represents the spearman correlation between TOBIAS/RNA. The TF clusters are  
 878 grouped into “Correlated” (Spearman $\geq$ 0.2), “shifted” (RNA max value appears before TOBIAS max value) and  
 879 “Not correlated” (Spearman $<$ 0.2 with no apparent shift in RNA).

880 **(b) Dynamic transcription factor binding during mouse embryonic development.** Similarly to figure 2A, the  
 881 heatmap depicts the TOBIAS-predicted footprint scores for 464 motifs during the time points 2C, 4C, 8C, ICM and  
 882 mESC. The rows are clustered into 6 clusters using hierarchical clustering. Individual cluster members are given  
 883 in Supplementary Table 3.



884

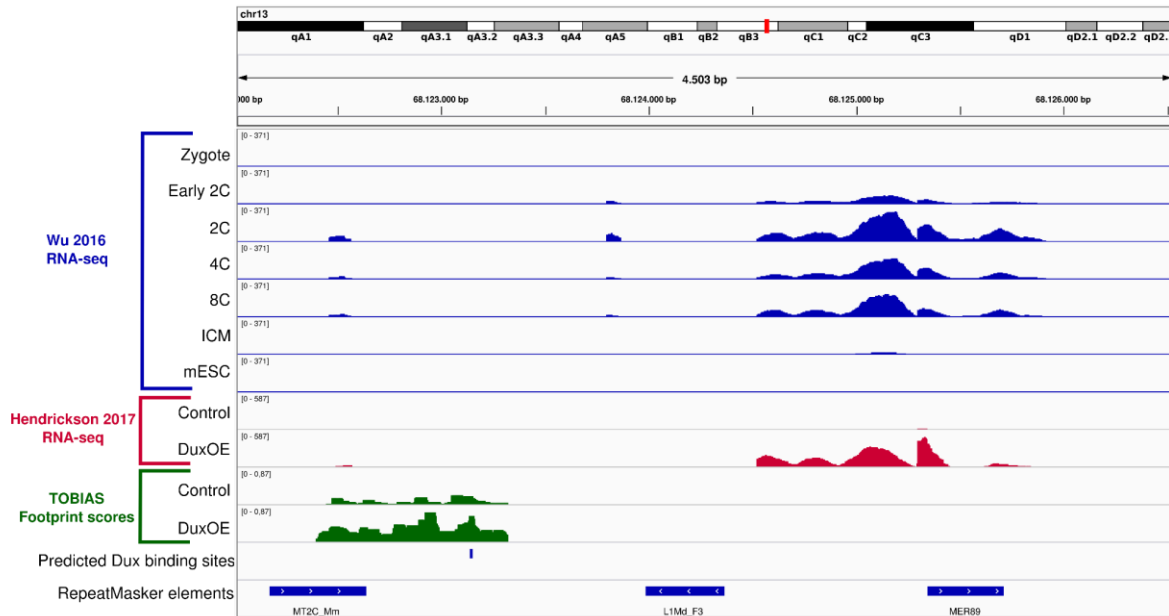
885 **Supplementary Figure 5: Dux binding is visible as footprints and correlate with ChIP-signal**

886 **(a) Correction of the Dux footprint using different bias correction methods.** The aggregate footprints for 12095  
 887 Dux binding sites (within ATAC-seq peaks) are shown between Control and DuxOE conditions. The top three  
 888 panels depict the uncorrected, expected and corrected signals as calculated by TOBIAS. The bottom panels depict  
 889 the same sites corrected by either HINT-ATAC or SeqOutBias methods.

890 **(b) A view of the footprinting scores in the promoter of Tdpoz1.** Genomic tracks show corrected ATAC-seq  
 891 cutsites at 1bp resolution (blue), footprint scores as calculated by TOBIAS (red), and pileup of reads from Dux  
 892 ChIP-seq of<sup>5</sup> (green). Potential Dux binding sites are highlighted in blue.

893 **(c-d) Footprinting correlates with ChIP-signal at multiple genomic loci.** Genomic tracks are the same as  
 894 described for (a).

895



896

897 **Supplementary Figure 6: Predicted Dux binding site correlates with increase in expression of close-**  
898 **by non-annotated regions**

899 *The figure shows genomic tracks of RNA-seq from <sup>21</sup> (blue) and <sup>5</sup> (red), TOBIAS footprint scores predicted from*  
900 *ATAC-seq (green) (<sup>5</sup>), predicted Dux binding site as well as known repeats as annotated by RepeatMasker (Smit,*  
901 *AFA, Hubley, R & Green, P. RepeatMasker Open-4.0).*

## 902 References

- 903 1. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J.  
904 Transposition of native chromatin for fast and sensitive epigenomic profiling of open  
905 chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**, 1213-  
906 1218 (2013).
- 907 2. Skene, P.J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution  
908 mapping of DNA binding sites. *eLife* **6**, e21856 (2017).
- 909 3. Eckersley-Maslin, M.A., Alda-Catalinas, C. & Reik, W. Dynamics of the epigenetic  
910 landscape during the maternal-to-zygotic transition. *Nat Rev Mol Cell Biol* **19**, 436-450  
911 (2018).
- 912 4. Jukam, D., Shariati, S.A.M. & Skotheim, J.M. Zygotic Genome Activation in  
913 Vertebrates. *Dev Cell* **42**, 316-332 (2017).
- 914 5. Hendrickson, P.G. et al. Conserved roles of mouse DUX and human DUX4 in  
915 activating cleavage-stage genes and MERVL/HERVL retrotransposons. *Nat Genet* **49**,  
916 925-934 (2017).
- 917 6. De Iaco, A. et al. DUX-family transcription factors regulate zygotic genome activation  
918 in placental mammals. *Nat Genet* **49**, 941-945 (2017).
- 919 7. Eckersley-Maslin, M.A. et al. MERVL/Zscan4 Network Activation Results in Transient  
920 Genome-wide DNA Demethylation of mESCs. *Cell Rep* **17**, 179-192 (2016).
- 921 8. Madisson, E. et al. Characterization and target genes of nine human PRD-like  
922 homeobox domain genes expressed exclusively in early embryos. *Sci Rep* **6**, 28995  
923 (2016).
- 924 9. Hesselberth, J.R. et al. Global mapping of protein-DNA interactions in vivo by digital  
925 genomic footprinting. *Nat Methods* **6**, 283-289 (2009).
- 926 10. Galas, D.J. & Schmitz, A. DNase footprinting: a simple method for the detection of  
927 protein-DNA binding specificity. *Nucleic Acids Res* **5**, 3157-3170 (1978).
- 928 11. Sung, M.H., Baek, S. & Hager, G.L. Genome-wide footprinting: ready for prime time?  
929 *Nat Methods* **13**, 222-228 (2016).
- 930 12. Vierstra, J. & Stamatoyannopoulos, J.A. Genomic footprinting. *Nat Methods* **13**, 213-  
931 221 (2016).
- 932 13. Karabacak Calviello, A., Hirsekorn, A., Wurmus, R., Yusuf, D. & Ohler, U. Reproducible  
933 inference of transcription factor footprints in ATAC-seq and DNase-seq datasets using  
934 protocol-specific bias modeling. *Genome Biology* **20**, 42 (2019).
- 935 14. Li, Z. et al. Identification of transcription factor binding sites using ATAC-seq. *Genome*  
936 *biology* **20**, 45-45 (2019).
- 937 15. Tripodi, I.J., Allen, M.A. & Dowell, R.D. Detecting Differential Transcription Factor  
938 Activity from ATAC-Seq Data. *Molecules* **23** (2018).
- 939 16. Baek, S., Goldstein, I. & Hager, G.L. Bivariate Genomic Footprinting Detects Changes  
940 in Transcription Factor Activity. *Cell Rep* **19**, 1710-1722 (2017).
- 941 17. Koster, J. & Rahmann, S. Snakemake-a scalable bioinformatics workflow engine.  
942 *Bioinformatics* **34**, 3600 (2018).
- 943 18. Di Tommaso, P. et al. Nextflow enables reproducible computational workflows. *Nat*  
944 *Biotechnol* **35**, 316-319 (2017).
- 945 19. Belmann, P. et al. de.NBI Cloud federation through ELIXIR AAI [version 1; peer review:  
946 2 approved, 1 not approved]. *F1000Research* **8** (2019).
- 947 20. Wu, J. et al. Chromatin analysis in human early development reveals epigenetic  
948 transition during ZGA. *Nature* **557**, 256-260 (2018).
- 949 21. Wu, J. et al. The landscape of accessible chromatin in mammalian preimplantation  
950 embryos. *Nature* **534**, 652-657 (2016).
- 951 22. Wang, K. & Nishida, H. REGULATOR: a database of metazoan transcription factors  
952 and maternal factors for developmental studies. *BMC Bioinformatics* **16**, 114 (2015).



- 953 23. Adjaye, J. & Monk, M. Transcription of homeobox-containing genes detected in cDNA  
954 libraries derived from human unfertilized oocytes and preimplantation embryos. *Mol*  
955 *Hum Reprod* **6**, 707-711 (2000).
- 956 24. Adhikary, S. et al. Miz1 is required for early embryonic development during  
957 gastrulation. *Mol Cell Biol* **23**, 7648-7657 (2003).
- 958 25. Home, P. et al. Genetic redundancy of GATA factors in the extraembryonic trophoblast  
959 lineage ensures the progression of preimplantation and postimplantation mammalian  
960 development. *Development* **144**, 876-888 (2017).
- 961 26. Xu, K. et al. Maternal Sall4 Is Indispensable for Epigenetic Maturation of Mouse  
962 Oocytes. *J Biol Chem* **292**, 1798-1807 (2017).
- 963 27. Svoboda, P. Mammalian zygotic genome activation. *Semin Cell Dev Biol* **84**, 118-126  
964 (2018).
- 965 28. Schulz, K.N. & Harrison, M.M. Mechanisms regulating zygotic genome activation.  
966 *Nature Reviews Genetics* **20**, 221-234 (2019).
- 967 29. Tohonen, V. et al. Novel PRD-like homeodomain transcription factors and  
968 retrotransposon elements in early human development. *Nat Commun* **6**, 8207 (2015).
- 969 30. Rhee, C. et al. ARID3A is required for mammalian placenta development.  
970 *Developmental Biology* **422**, 83-91 (2017).
- 971 31. Winger, Q., Huang, J., Auman, H.J., Lewandoski, M. & Williams, T. Analysis of  
972 Transcription Factor AP-2 Expression and Function During Mouse Preimplantation  
973 Development1. *Biology of Reproduction* **75**, 324-333 (2006).
- 974 32. Pastor, W.A. et al. TFAP2C regulates transcription in human naive pluripotency by  
975 opening enhancers. *Nat Cell Biol* **20**, 553-564 (2018).
- 976 33. Eckersley-Maslin, M. et al. Dppa2 and Dppa4 directly regulate the Dux-driven zygotic  
977 transcriptional program. *Genes Dev* **33**, 194-208 (2019).
- 978 34. De Iaco, A., Coudray, A., Duc, J. & Trono, D. DPPA2 and DPPA4 are necessary to  
979 establish a 2C-like state in mouse embryonic stem cells. *EMBO Rep* **20** (2019).
- 980 35. Whiddon, J.L., Langford, A.T., Wong, C.J., Zhong, J.W. & Tapscott, S.J. Conservation  
981 and innovation in the DUX4-family gene network. *Nat Genet* **49**, 935-940 (2017).
- 982 36. Huang, C.J., Chen, C.Y., Chen, H.H., Tsai, S.F. & Choo, K.B. TDPOZ, a family of  
983 bipartite animal and plant proteins that contain the TRAF (TD) and POZ/BTB domains.  
984 *Gene* **324**, 117-127 (2004).
- 985 37. Berest, I. et al. Quantification of differential transcription factor activity and multiomics-  
986 based classification into activators and repressors: *&lt;em>diffTF&lt;/em>*  
987 *bioRxiv*, 368498 (2018).
- 988 38. Lee, S.-E., Lee, S.-Y. & Lee, K.-A. RhoX in mammalian reproduction and development.  
989 *Clin Exp Reprod Med* **40**, 107-114 (2013).
- 990 39. Borgmann, J. et al. The human RHOX gene cluster: target genes and functional  
991 analysis of gene variants in infertile men. *Hum Mol Genet* **25**, 4898-4910 (2016).
- 992 40. Royall, A.H., Maeso, I., Dunwell, T.L. & Holland, P.W.H. Mouse Obox and Crxos  
993 modulate preimplantation transcriptional profiles revealing similarity between  
994 paralogous mouse and human homeobox genes. *Evodevo* **9**, 2 (2018).
- 995 41. Percharde, M. et al. A LINE1-Nucleolin Partnership Regulates Early Development and  
996 ESC Identity. *Cell* **174**, 391-405.e319 (2018).
- 997 42. Buenrostro, J.D. et al. Single-cell chromatin accessibility reveals principles of  
998 regulatory variation. *Nature* **523**, 486-490 (2015).
- 999 43. Harrison, P.W. et al. The European Nucleotide Archive in 2018. *Nucleic Acids Res*  
1000 (2018).
- 1001 44. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing  
1002 reads. *2011* **17**, 3 (2011).
- 1003 45. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21  
1004 (2013).
- 1005 46. Zerbino, D.R. et al. Ensembl 2018. *Nucleic Acids Res* **46**, D754-D761 (2018).
- 1006 47. Feng, J.X., Liu, T., Qin, B., Zhang, Y. & Liu, X.S. Identifying ChIP-seq enrichment using  
1007 MACS. *Nat Protoc* **7**, 1728-1740 (2012).

- 1008 48. Frankish, A. et al. GENCODE reference annotation for the human and mouse  
1009 genomes. *Nucleic Acids Research* **47**, D766-D773 (2018).
- 1010 49. Kondili, M. et al. UROPA: a tool for Universal ROBust Peak Annotation. *Sci Rep* **7**,  
1011 2593 (2017).
- 1012 50. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion  
1013 for RNA-seq data with DESeq2. *Genome biology* **15**, 550-550 (2014).
- 1014 51. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina  
1015 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 1016 52. Khan, A. et al. JASPAR 2018: update of the open-access database of transcription  
1017 factor binding profiles and its web framework. *Nucleic Acids Res* **46**, D1284 (2018).
- 1018 53. Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription  
1019 factor binding models for human and mouse via large-scale ChIP-Seq analysis.  
1020 *Nucleic Acids Research* **46**, D252-D259 (2017).
- 1021 54. Sebastian, A. & Contreras-Moreira, B. footprintDB: a database of transcription factors  
1022 with annotated cis elements and binding interfaces. *Bioinformatics* **30**, 258-265 (2013).
- 1023 55. Yin, Y. et al. Impact of cytosine methylation on DNA binding specificities of human  
1024 transcription factors. *Science* **356**, eaaj2239 (2017).
- 1025 56. Machanick, P. & Bailey, T.L. MEME-ChIP: motif analysis of large DNA datasets.  
1026 *Bioinformatics* **27**, 1696-1697 (2011).
- 1027 57. Durinck, S., Spellman, P.T., Birney, E. & Huber, W. Mapping identifiers for the  
1028 integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat Protoc*  
1029 **4**, 1184 (2009).
- 1030 58. Gu, L. et al. The Histone Demethylase PHF8 Is Essential for Endothelial Cell Migration.  
1031 *PLoS One* **11**, e0146645 (2016).
- 1032 59. Robinson, J.T. et al. Integrative genomics viewer. *Nature biotechnology* **29**, 24-26  
1033 (2011).
- 1034 60. Egorov, A.A. et al. svist4get: a simple visualization tool for genomic tracks from  
1035 sequencing experiments. *BMC Bioinformatics* **20**, 113 (2019).
- 1036 61. Shannon, P. et al. Cytoscape: A Software Environment for Integrated Models of  
1037 Biomolecular Interaction Networks. *Genome Research* **13**, 2498-2504 (2003).
- 1038 62. Ramírez, F. et al. deepTools2: a next generation web server for deep-sequencing data  
1039 analysis. *Nucleic acids research* **44**, W160-W165 (2016).
- 1040 63. Siddharthan, R. Dinucleotide weight matrices for predicting transcription factor binding  
1041 sites: generalizing the position weight matrix. *PLoS One* **5**, e9722 (2010).
- 1042 64. Koohy, H., Down, T.A. & Hubbard, T.J. Chromatin accessibility data sets show bias  
1043 due to sequence specificity of the DNase I enzyme. *PLoS One* **8**, e69853 (2013).
- 1044 65. Gusmao, E.G., Allhoff, M., Zenke, M. & Costa, I.G. Analysis of computational  
1045 footprinting methods for DNase sequencing experiments. *Nat Methods* **13**, 303-309  
1046 (2016).
- 1047 66. Korhonen, J.H., Palin, K., Taipale, J. & Ukkonen, E. Fast motif matching revisited: high-  
1048 order PWMs, SNPs and indels. *Bioinformatics* **33**, 514-521 (2016).
- 1049 67. Benaglia, T., Chauveau, D., Hunter, D.R. & Young, D.S. mixtools: An R Package for  
1050 Analyzing Mixture Models. *Journal of Statistical Software; Vol 1, Issue 6 (2010)* (2009).
- 1051 68. Neph, S. et al. An expansive human regulatory lexicon encoded in transcription factor  
1052 footprints. *Nature* **489**, 83-90 (2012).
- 1053 69. Sung, M.H., Guertin, M.J., Baek, S. & Hager, G.L. DNase footprint signatures are  
1054 dictated by factor dynamics and DNA sequence. *Mol Cell* **56**, 275-285 (2014).
- 1055 70. Piper, J. et al. Wellington: a novel method for the accurate identification of digital  
1056 genomic footprints from DNase-seq data. *Nucleic Acids Res* **41**, e201 (2013).
- 1057 71. Boyle, A.P. et al. High-resolution genome-wide in vivo footprinting of diverse  
1058 transcription factors in human cells. *Genome Res* **21**, 456-464 (2011).
- 1059 72. Pique-Regi, R. et al. Accurate inference of transcription factor binding from DNA  
1060 sequence and chromatin accessibility data. *Genome Res* **21**, 447-455 (2011).
- 1061 73. Luo, K. & Hartemink, A.J. Using DNase digestion data to accurately identify  
1062 transcription factor binding sites. *Pac Symp Biocomput*, 80-91 (2013).

- 1063 74. Kahara, J. & Lahdesmaki, H. BinDNase: a discriminatory approach for transcription  
1064 factor binding prediction using DNase I hypersensitivity data. *Bioinformatics* **31**, 2852-  
1065 2859 (2015).
- 1066 75. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nature*  
1067 *methods* **9**, 357-359 (2012).
- 1068 76. Bailey, T.L. & Machanick, P. Inferring direct DNA binding from ChIP-seq. *Nucleic Acids*  
1069 *Research* **40**, e128-e128 (2012).
- 1070 77. Martins, A.L., Walavalkar, N.M., Anderson, W.D., Zang, C. & Guertin, M.J. Universal  
1071 correction of enzymatic sequence bias reveals molecular signatures of protein/DNA  
1072 interactions. *Nucleic Acids Res* **46**, e9 (2018).
- 1073 78. Raj, A., Shim, H., Gilad, Y., Pritchard, J.K. & Stephens, M. msCentipede: Modeling  
1074 Heterogeneity across Genomic Sites and Replicates Improves Accuracy in the  
1075 Inference of Transcription Factor Binding. *PLoS One* **10**, e0138030 (2015).
- 1076 79. Sherwood, R.I. et al. Discovery of directional and nondirectional pioneer transcription  
1077 factors by modeling DNase profile magnitude and shape. *Nature biotechnology* **32**,  
1078 171-178 (2014).  
1079