

Widespread conservation of chromatin accessibility patterns and transcription factor binding in human and chimpanzee induced pluripotent stem cells

Irene Gallego Romero^{1,2*}, Shyam Gopalakrishnan³, Yoav Gilad^{4,5*}

1. Melbourne Integrative Genomics,
University of Melbourne
Melbourne, Australia,

2. School of BioSciences,
University of Melbourne
Melbourne, Australia,

3. Natural History Museum of Denmark,
University of Copenhagen
Copenhagen, Denmark

4. Department of Medicine,
University of Chicago
Chicago, USA

5. Department of Human Genetics,
University of Chicago
Chicago, USA

* To whom correspondence should be addressed: irene.gallego@unimelb.edu.au,
gilad@uchicago.edu

1 **Abstract:**

2 Changes in gene regulation have been shown to contribute to phenotypic differences between
3 closely related species, most notably in primates. It is likely that a subset of inter-species
4 regulatory differences can be explained by changes in chromatin accessibility and transcription
5 factor binding, yet there is a paucity of comparative data sets with which to investigate this.
6 Using ATAC-seq, we profiled genome-wide chromatin accessibility in a matched set of 6 human
7 and 6 chimpanzee (*Pan troglodytes*, our closest living relative) induced pluripotent stem cells
8 from which we have previously collected gene expression data. We examined chromatin
9 accessibility patterns near 20,745 orthologous transcription start sites and used a footprinting
10 algorithm to predict transcription factor binding activity in each species. We found that the
11 majority of chromatin accessibility patterns and transcription factor activity are conserved
12 between these two closely related species. Interestingly, interspecies divergence in chromatin
13 accessibility and transcription factor binding in pluripotent cells appear to contribute not to
14 differences in the pluripotent state, but to downstream developmental processes. Put together,
15 our findings suggest that the pluripotent state is extremely stable and potentially subject to
16 stronger evolutionary constraint than other somatic tissues.

17

18

19 **Abbreviations:**

20 ATAC: assay for transposase-accessible chromatin
21 CPM: counts per million
22 DA: differentially accessible
23 DE: differentially expressed
24 ESC: embryonic stem cell
25 FDR: false discovery rate
26 iPSC: induced pluripotent stem cell
27 MPBS: motif-predicted binding site(s)
28 PCA: principal component analysis
29 PWM: position weight matrix
30 RPKM: reads per kilobase per million reads mapped
31 TF: transcription factor

32 TFBS: transcription factor binding site(s)

33 TSS: transcription start site(s)

34

35

36 **Keywords:**

37 Transcription factor, chromatin, chimpanzee, human, induced pluripotent stem cells, footprinting

38

39

40 **Introduction:**

41 The contributions of gene regulatory differences to speciation and adaptation have been
42 documented time and again by comparative genomics, aligning with the long-standing intuition
43 that turnover at protein coding regions might not be sufficient to account for much of the
44 phenotypic change among closely related species [1-5]. For example, ultra-conserved non-coding
45 regions that exhibit accelerated change in a single lineage [6-9] have been shown to contribute to
46 morphological differences between species [10-13]. In cases where the mechanism of action has
47 been studied, the gain or loss of a single regulatory element can often be sufficient to drive
48 phenotypes as substantial as the loss of limb development in snakes [14], or the loss of armour
49 plates in fish [15].

50 Humans and other primates are no exception [5]. Human-specific gene expression and
51 regulatory differences have been characterised during brain [16, 17] and limb development [18,
52 19]; likewise, steady-state differences in gene expression levels across myriad tissues have been
53 repeatedly identified between humans and other primates [20-25] and attributed to changes in
54 regulatory mechanisms such as histone modifications [26, 27], DNA methylation [28-30], repeat
55 expansions [31, 32] (but see [33]), and differences in the activity of specific transcription factors
56 (TF) [34]. A comparative study of DNaseI hypersensitive sites found more than 1000 loci that
57 showed either human-specific or chimpanzee-specific gains or losses in hypersensitivity in
58 fibroblasts and immortalised lymphoblastoid cells [35]. Many of these sites, which mark open
59 chromatin presumed to be accessible to the cell's transcriptional machinery, were located close to
60 genes that also exhibited interspecies differences in their expression levels. Put together, a strong
61 body of evidence suggests that regulatory changes and the associated differences in chromatin
62 state make functional contributions to phenotypic differences between these species.

63 Two of the most intuitive mechanisms by which interspecies differences in gene
64 expression levels can arise are through changes in chromatin accessibility and TF binding
65 activity. Regulatory networks exhibit high degrees of robustness and resilience such that gene
66 expression levels are generally conserved across species [36, 37], but multiple studies have
67 shown that turnover at individual transcription factor binding sites (TFBS) is rapid, with few
68 binding sites conserved across large evolutionary distances [38]. For example, a mouse model of
69 human trisomy 21 carrying the entire human chromosome 21 alongside two copies of the mouse
70 orthologue showed that expression levels of the genes on the human chromosome were more

71 similar to those seen in human cells than to expression levels of their murine orthologues in the
72 same cell. This finding argues strongly for a regulatory environment that is driven by DNA
73 sequence rather than the nuclear environment [39]. Similarly, research in *Drosophila* has shown
74 that changes in chromatin accessibility can be associated with changes in TF binding activity
75 [40], and that chromatin accessibility and of TF binding sites can be used to predict gene
76 expression divergence [41].

77 However, these previous studies have focused on only a handful of TFs, chosen because
78 of their functional importance, rather than perform a general survey of the regulatory landscape.
79 Direct high-throughput measurements of binding activity (using ChIP-seq, for example) across a
80 representative suite of TFs in multiple samples – be they drawn from different tissues,
81 individuals, or species – remain rare. In spite of recent technological advances, the need for
82 validated antibodies, and/or large amounts of starting material often makes generating these data
83 sets prohibitively expensive both in terms of labour and cost. An alternative popular approach,
84 footprinting, combines chromatin accessibility data (originally in the form of DNaseI cleavage)
85 with knowledge of a TF's preferred binding sequence to predict whether a given site in the
86 genome is bound or not [42].

87 Here, we profiled chromatin accessibility using ATAC-seq [43] in 6 human and 6
88 chimpanzee iPSC lines, which we had previously generated and extensively validated [44]. We
89 first used these data to examine interspecies patterns of chromatin accessibility at 20,745
90 orthologous transcription start sites, and we examined the correlation of these data with existing
91 gene expression data from the same cell lines. Then we applied a TF footprinting algorithm [45,
92 46] to 306 known position weight matrices (PWM), in order to perform an unprecedented
93 characterisation of binding activity at over 130 million putative TF binding sites genome-wide.
94 We found that the majority of chromatin accessibility and TF binding activity is conserved
95 between these two closely related species, but in instances where it is not, interspecies change
96 appears to contribute to downstream developmental processes.

97

98 **Results:**

99 To characterise the chromatin landscape in chimpanzee and human iPSCs, we generated
100 ATAC-seq [43] libraries from six previously described lines from each species [44], and
101 sequenced each sample to an average depth of ~252 million reads/sample. After exhaustive

102 quality control and filtering of duplicate reads, mtDNA reads and reads outside high confidence
103 orthologous regions (see methods), we retained an average of ~19.7 million reads per sample that
104 mapped to a set of defined high-quality orthologous regions in the chimpanzee and human
105 genomes (see methods and supplementary table 1). We additionally used the average ~99 million
106 mtDNA reads from each sample to reconstruct their full mitochondrial genome sequences and
107 infer the maternal lineage of our chimpanzee samples (see methods and supplementary figure 1),
108 which was previously unknown.

109

110 *Changes in chromatin accessibility and gene expression levels*

111 We initially focused our analyses on genomic regions that are expected to be accessible:
112 transcription start sites. We began by defining 20,745 2kb windows centred on orthologous
113 transcription start sites (orthologous TSS; see methods for more details). We also defined an
114 equivalent set of control orthologous background genomic regions matched for broad
115 mappability and located at least 5kb from any known TSS. Principal component analysis (PCA)
116 of the orthologous TSS dataset does not show a clear visual separation between species
117 (supplementary figure 2). While species is not significantly associated with any PC using the
118 chromatin accessibility data from the orthologous background regions, only PC2 is significantly
119 associated with species in the orthologous TSS dataset (PC2 $P = 0.007$; see supplementary table
120 2 for associations with additional PCs).

121 The results of the statistical analysis and the lack of clear visual separation of species in
122 our data ran counter to previous observations of strong separation between species along the first
123 PC when other types of regulatory data are used [44]. We thus sought an explanation. We
124 compared the distribution of CPM values at orthologous TSS and orthologous background
125 regions, and found that the former is clearly, and unexpectedly, bimodal (figure 1a;
126 supplementary figure 3). By using k-means clustering on each species' mean \log_2 CPM values,
127 we found that of the 20,745 orthologous TSS in our dataset, 5,675 were assigned to a secondary
128 CPM peak (figure 1b), which is associated with highly accessible regions. PCA using the data
129 from only these 5,675 highly accessible TSS reveals a much stronger association with species
130 than before (PC1 $P = 4.0 \times 10^{-5}$; figure 1c), which can also be easily ascertained visually. The lack
131 of clear separation of species in our overall data, therefore, has a technical explanation related to

132 data quality (see methods for more details and additional relevant QC analysis), which impacts
133 certain downstream analyses as we discuss below.

134 To consider the ATAC-seq data alongside corresponding gene expression measurements,
135 we re-analysed previously published RNA-sequencing data from the same cell lines [44], and
136 mapped reads to an updated reference orthologous transcriptome (see methods and [21]). We
137 detected the expression of 12,674 genes at \log_2 CPM ≥ 1 in at least half the individuals from
138 one species, yet we were able to confidently identify an orthologous TSS for only 4,210 (33%) of
139 these genes. This is due to our stringent definition of orthologous TSS, but the alternative is to
140 introduce a strong bias towards regions with better mappability in humans than chimpanzees,
141 which would impact our ability to detect genuine interspecies differences in chromatin
142 accessibility. Of the expressed genes for which we can identify an orthologous TSS with
143 confidence, 3,150 (75.0%) were genes whose orthologous TSS was part of the 'highly accessible'
144 group defined above, a significant excess (hypergeometric $P < 10^{-16}$) that supports this subset as
145 being enriched for true biological signal (figure 1b).

146 Levels of accessibility at an orthologous TSS and the expression levels of the
147 corresponding gene within an individual are weakly correlated (with Spearman's ρ ranging
148 between 0.05 for line Hutt60 to 0.17 for line H20961; $P < 7 \times 10^{-4}$ for all individuals). The
149 association improves somewhat if we only consider the 3,150 expressed genes with an
150 orthologous TSS within the 'highly accessible' subset, but still remains modest: ρ ranges from
151 0.08 (line Hutt60; $P = 2.6 \times 10^{-6}$) to 0.25 (line C8861G; $P = 1.7 \times 10^{-46}$). This observation suggests
152 that chromatin accessibility acts in a somewhat coarse fashion – remodelling is necessary to
153 allow for transcription to occur, but the fine-tuning of expression levels likely occurs through
154 additional mechanisms.

155 Bearing this in mind, we tested the 5,675 'highly accessible' orthologous TSS (regardless of
156 whether these TSS were associated with an expressed gene) for inter-species differences in
157 accessibility, using an identical framework to that which we have previously used to detect
158 differential gene expression [44]. At an FDR of 5%, 1,598 orthologous TSS are differentially
159 accessible (DA) between the two species (supplementary table 3). Of these, 313 differentially
160 accessible TSS are also associated with differentially expressed genes (supplementary table 4), a
161 significant overlap (hypergeometric $P = 8.6 \times 10^{-3}$, figure 2a). However, the relationship between
162 inter-species differential expression and differential accessibility is not straightforward: Only in

163 64.5% of cases the inter-species direction of the effect is as expected, with the effect size sign
164 being the same in in both the chromatin accessibility and gene expression data sets (figure 2b).

165 The 1,598 genes associated with orthologous TSS that are differentially accessible
166 between species are enriched for five Gene Ontology Biological Process [47] terms, GO0007275:
167 multicellular organismal development, GO0048731: system development, GO0044767: single-
168 organism developmental process, GO0032502: developmental process and GO0060429:
169 epithelium development (FDR = 10%). These enrichments are driven by a largely overlapping
170 set of 398 genes, 294 of which are detectably expressed in at least one species, and 90 of which
171 are differentially expressed between the two species (figure 2c; a full list is provided as
172 supplementary table 5).

173 Interestingly, 104 of the genes associated with the orthologous TSS driving our GO results
174 are not detected as expressed in iPSCs from either species. In order to understand the factors
175 driving this observation, we asked whether these 104 regions might be playing a role in
176 establishing gene expression patterns following exit from the pluripotent state [48]. To do so, we
177 considered our data in combination with information on bivalent chromatin regions
178 ("10_TssBiv", "11_BivFlkn" and "12_EnhBiv") identified in human embryonic stem cell line H1
179 as part of the Roadmap Epigenomics Project [49]. We found that 83 of these 104 orthologous
180 TSS overlap at least one bivalent region in the Roadmap Epigenomics data, a far higher fraction
181 than expected by chance alone (hypergeometric $P < 10^{-26}$), as well as a high fraction compared to
182 the overlap of bivalent regions with orthologous TSS associated with genes expressed in at least
183 one species (75 out of 292; $P = 0.97$) or with TSS associated with differentially expressed genes
184 (30 out of 90; $P = 0.26$). Indeed, the 104 genes associated with these TSS include multiple TFs
185 with well-established roles in early development such as *ISLI*, which plays key roles in cardiac
186 development [50], or *NEUROD2* and *NEUROD4*, both implicated in development of neural
187 tissue [51-53], as well as other genes of less clear role, such as *RTN4RL1*.

188

189 ***Using ATAC-seq to infer inter-species differences in transcription factor binding activity***

190 Because chromatin accessibility is not highly correlated with gene expression patterns
191 either within or between species, we sought to characterise the landscape of TF binding activity
192 in these cells by using footprinting analysis. To do so, we used a recent extension of the
193 CENTIPEDE algorithm (msCentipede) that can better account for signal heterogeneity across

194 sites (see methods; [46]), and predicted binding activity genome-wide across 306 position weight
195 matrices (PWMs) from the HOCOMOCO database [54]. We considered 133,103,977 motif-
196 predicted binding sites (MPBS) with a PWM score ≥ 7 in both the human (hg19) and
197 chimpanzee (panTro2.1.3) genomes. This metric, which we calculated for each MPBS in each
198 species, reflects how closely the locus matches the ideal PWM motif, with higher scores
199 denoting higher fidelity. While each PWM is nominally associated with a single TF, there can be
200 redundancy between PWMs associated with closely related TFs; thus we refer primarily to
201 PWMs rather than TFs in the remaining sections.

202 We began by assessing the validity of our approach by comparing our binding predictions
203 with all publicly released ENCODE [55] TF ChIP-seq data from the human embryonic stem cell
204 line H1 ($n = 28$) and the human induced pluripotent stem cell line GM23338 ($n = 5$). Although
205 recall and precision vary widely for different PWMs, we found that binding predictions from
206 msCentipede often recapitulate findings for well-characterised PWMs and TFs (figure 3;
207 additional details are available in supplementary table 6). For instance, ENCODE data contains
208 three independent CTCF ChIP-seq experiments, one performed on H1 hESCs and two on
209 GM23338 hiPSCs. On average, 25.1% of genome-wide CTCF MPBS in our data overlap
210 ENCODE CTCF ChIP-seq peaks across at least one of the three experiments. When we
211 considered only those sites msCentipede predicts to be bound in our human iPSCs, 73.7%
212 overlap an ENCODE CTCF ChIP-seq peak, and 98.3% when we only considered bound MPBS
213 with a PWM score ≥ 12 (a widely-accepted threshold for high quality MPBS). We found similar
214 concordance when we considered MPBS for REST (also known as NRSF), NRF1, YY1 and
215 GABPA amongst others.

216 Having validated our approach, we turned our attention to a comparison of MPBS across
217 species. Across all PWMs, an average of 16.9% of sites are classified as bound (see methods) in
218 humans, although values range from 0.7% to 97.3% for specific PWMs. In chimpanzee, an
219 average of 9.5% of sites per PWM are classified as bound, ranging from 0.4% to 98.3%. The
220 inter-species overall difference in binding can be explained by the lower read coverage in our
221 chimpanzee samples, as discussed above (and in more detail in the methods). Reassuringly, an
222 average of 85.5% of sites predicted to be bound in chimpanzees are also bound in humans (figure
223 4a). Thus we interpret the bulk of observed human-unique binding events as not being

224 biologically meaningful but rather driven by technical differences, although we discuss some
225 exceptions below.

226 We limited our next analysis to the subset of motifs with a PWM score ≥ 12 ($n = 906,612$
227 across 173 PWMs). We found that, on average, 23.3% of sites are classified as bound in humans
228 and 13.0% in chimpanzees. Of sites bound in chimpanzees, of 86.6% are also classified as bound
229 in humans (figure 4b). We observed that predicted binding events do not occur randomly
230 throughout the genome. When we classified binding sites by their distances to the nearest
231 orthologous TSS, we found that binding significantly increases in frequency as distance to an
232 orthologous TSS decreases ($P < 1*10^{-16}$; figure 4c). The concordance between human and
233 chimpanzee results also increases for binding sites that are closer to the TSS. Interestingly, this
234 observation is consistent across the vast majority of PWMs; the only exceptions are the PWMs
235 associated with ERF and OTX2 in both species, PBX1 in humans only, and HINFP, TCF7L2 and
236 NR3C2 in chimpanzees only (supplementary figure 4). Of these six TFs, *HINFP* and *TCF7L2*
237 are differentially expressed between humans and chimpanzees; *ERF* is differentially accessible
238 between the two species. However, all six PWMs are associated with small numbers of high
239 quality MPBS (under 600 in all cases, as low as 32 in the case of HINFP), which may impact the
240 predictive accuracy of msCentipede in these cases.

241

242 *Unique binding patterns across species*

243 There are multiple instances of predicted species-specific TF binding activity in our data.
244 We explored possible mechanisms for divergence in TF binding by considering four scenarios: 1.
245 *trans*-acting interspecies differences in the TF itself, possibly indicative of a change in motif
246 preference, summarised by dN/dS values; 2. differences in the expression levels of the TF that
247 binds the motif, also suggestive of *trans*-acting change; 3. differences in chromatin accessibility
248 at MPBS, which can be tested by asking whether there is an excess of chimpanzee-unique
249 binding events in regions differentially accessible between the species; and 4. *cis*-acting
250 sequence turnover in the MPBS. To test this last possibility, we defined a simple metric, Δ PWM
251 score, which is the difference at each MPBS between the PWM score in humans and
252 chimpanzees.

253 Neither dN/dS ratios of TF divergence nor inter-species differences in the expression
254 level of a TF are significantly associated with species-specific binding genome-wide. In contrast

255 Δ PWM score between human and chimpanzee is associated with inter-species differences in TF
256 binding, broadly suggesting that many predicted binding differences are due to *cis*-acting
257 mechanisms. Sites inferred to be bound only in chimpanzees have, on average, a higher PWM
258 score in chimpanzees than in humans, and vice versa (figure 5a). While the distribution of the
259 Δ PWM score metric for sites bound in both species is centred at 0.002 ($n = 146,554$), the Δ PWM
260 distribution for sites bound only in chimpanzee or only in human are clearly skewed (mean
261 Δ PWM score for sites bound only in humans = 0.25; $n = 107,784$; $P < 10^{-16}$; for sites bound only
262 in chimpanzees = -0.29; $n = 14,803$; $P < 10^{-16}$).

263 We then focused on chimpanzee-specific binding events, as we expect these to more
264 likely be true positive unique events (in contrast to the human-specific events, many of which are
265 likely to be missing from chimpanzee due to technical reasons, as previously discussed). We
266 asked whether chimpanzee-specific binding events around orthologous TSSs can be associated
267 with differences in corresponding gene expression levels or chromatin accessibility. We
268 considered all orthologous TSS (± 5 kb) with a least one predicted high-quality chimpanzee-
269 binding event (but not necessarily chimpanzee-specific) and with matched accessibility ($n =$
270 3,870) or expression ($n = 2,981$) data. We found a weak but significant correlation (Pearson's $R =$
271 0.06; $P = 2.0 \times 10^{-4}$) between interspecies differences in accessibility around orthologous TSS and
272 the number of chimpanzee-specific binding events in the region. The magnitude of these
273 differences is small: the mean number of chimpanzee-specific binding events for orthologous
274 TSS that are not differentially accessible (DA) is 0.57 ($n = 3,326$), 0.83 for all DA orthologous
275 TSS ($n = 544$), and 1.42 for DA orthologous TSS with an absolute \log_2 fold change in
276 accessibility ≥ 2 ($n = 19$). We also found a weak correlation (Pearson's $R = 0.05$, $P = 0.01$)
277 between interspecies gene expression effect size (absolute \log_2 fold change) and the number of
278 chimpanzee-specific binding events occurring within 5kb of the corresponding orthologous TSS.
279 Again, the size of this effect is small: the mean number of chimpanzee-specific binding events
280 for non-DE genes is 0.44 ($n = 2,043$), 0.53 for all DE genes ($n = 938$), and 1.03 for DE genes
281 with an absolute \log_2 fold change ≥ 2 ($n = 62$).

282 We asked whether any PWM is associated with a systematic excess or deficit of species-
283 specific binding genome-wide, which could suggest broad regulatory network rewiring and
284 turnover. Again we focused on chimpanzee-specific excesses, and considered 109 PWMs with at
285 least 100 predicted binding events in chimpanzees at high quality MPBS. We found 6 PWMs

286 that exhibit chimpanzee-specific binding at levels more than 2 standard deviations away from the
287 mean, associated with the TFs BARX1, CPEB1, FOXO1, PAX7, POU6F1 and ZFH3 (figure
288 5b). These PWMs are ranked as class D by HOCOMOCO, suggestive of poor predictive value,
289 with the exception of FOXO1, which is class C (see methods for more details). We did not find
290 any PWMs that exhibit excessive human-specific binding. However, if we simply rank PWMs
291 by their amount of human-specific binding we found that, with the exception of FOXO1, the
292 same PWMs as above are amongst those with the greatest amount of specific binding in humans
293 (figure 5b) potentially suggesting that these PWM generally evolve rapidly in the two species, or,
294 alternatively, that in these cases the PWM does not capture the full scope of binding activity.

295 Of the 6 TFs associated with these PWMs, 3 are differentially expressed between human
296 and chimpanzee iPSCs – *FOXO1*, *BARX1* and *PAX7*, although we found no excess of
297 differentially expressed genes relative to the dataset-wide average (permutation $P = 0.7$). This
298 observation might be explained by high redundancy in the binding motifs. For example, 399 of
299 1,014 binding events associated with the FOXO1 PWM in chimpanzee iPSCs are unique to the
300 species. But the PWM itself has a high similarity to other sites associated with different forkhead
301 family members such as *FOXJ3*, which is expressed in both humans and chimpanzees at high
302 levels and with which *FOXO1* shares ~30% of possible high-quality binding sites.

303

304 *Dynamics of activity in the core pluripotency regulatory network*

305 Finally, we focused on PWMs associated with TFs with key roles in maintaining the
306 pluripotent state. The TFs *OCT4* (also known as *POU5F1*), *SOX2* and *NANOG* sit atop the gene
307 regulatory network of pluripotency, and act cooperatively to maintain this state [56]. Although
308 both our previous work and that of others [44, 57] has shown that the pluripotency gene
309 regulatory network is highly conserved between humans and chimpanzees, we found that
310 conservation in the binding locations of these PWMs is not particularly high. Only 57.9% of
311 *OCT4*, 55.1% of *SOX2* and 55.2% of *NANOG* sites in the high-quality motif subset predicted to
312 be bound in chimpanzee are also predicted to be bound in human. These observations, however,
313 must be tempered by the poor concordance between the sites msCentipede predicts to be bound
314 and ChIP-seq peaks for *OCT4* and *NANOG* (both profiled in GM23338 iPSCs) in ENCODE
315 data (figure 3), which in turn are likely driven by the complex binding dynamics of these TFs
316 [58].

317 We also considered a larger set of 15 master pluripotency regulators drawn from the
318 literature (specifically, from [56, 59, 60]). In this case, we found only few instances of
319 chimpanzee-specific binding, suggesting overall high conservation of pluripotency pathways
320 between the two species (mean fraction of chimpanzee-specific binding = 4.6% at high-quality
321 sites, 4.9% across all sites; figure 6). This result is consistent with our previous observation that
322 master pluripotency regulators are generally not differentially expressed between iPSCs from the
323 two species [44].

324

325 **Discussion:**

326 The regulation of gene expression is a dynamic process orchestrated by multiple
327 interacting cellular mechanisms. Transcription factors operate at the point of transcriptional
328 initiation, by promoting (or sometimes inhibiting) the formation of the transcriptional complex
329 through a suite of mechanisms. Here we have characterised and inferred the chromatin and TF
330 binding landscapes in human and chimpanzee iPSCs. As might be expected given that we chose
331 to use a cell type that mimics very early embryonic development, overall we observed high
332 conservation of chromatin accessibility and TF binding in the two species. We found relatively
333 few instances of interspecies regulatory differences that can be associated to downstream
334 differences in expression. Taken together, our results consistently suggest that the pluripotency
335 gene regulatory network is highly conserved between humans and chimpanzees – and seems
336 highly robust to those interspecies differences we do observe. Indeed, there was little evidence to
337 suggest that differences in chromatin architecture impact the pluripotency gene regulatory
338 network that helps maintain cell identity. Much higher regulatory divergence is typically
339 observed in differentiated tissues or cell types [21, 34] and indeed, we have also observed inter-
340 species differences in chromatin accessibility near genes implicated in early developmental
341 processes, such as embryonic patterning. The systematic difference in chromatin accessibility
342 amongst bivalently modified genes known to be implicated in embryonic patterning is intriguing
343 and, in our mind, worthy of further investigation.

344 We performed an in-depth examination of the potential of footprinting to infer genome-
345 wide TF binding activity across scores of TFs simultaneously, using publicly available PWM
346 datasets. There are variable levels of concordance between our binding predictions and
347 ENCODE ChIP-seq data. While some of this discordance is attributable to technical artefacts,

348 our data are in line with previous observations that PWMs do not always fully capture the
349 binding behaviours of particular TFs [61]. In spite of this caveat, our results suggest that much of
350 the binding activity is conserved between human and chimpanzees across the vast majority of
351 PWMs we tested. This finding recapitulates a study of early development between the closely
352 related *Drosophila melanogaster* and *Drosophila yakuba*, which found that these species shared
353 between 85–98% of TFBS across six developmental regulators [62].

354 When considered at the scale of entire gene regulatory networks, we found strong
355 evidence of redundancy in TF binding activity [36]. We recapitulated the observation that *cis*-
356 acting turnover at possible binding sites is the main driver of divergence in binding activity [38],
357 more than other possible forces such as *trans*-acting evolutionary divergence at the TF itself, or
358 differences in TF expression levels. Indeed, the majority of the differences we identified, both in
359 chromatin accessibility and TF binding, appear to be *cis*-acting. It has been proposed that *cis*-
360 acting changes are more likely than *trans*-acting changes to drive evolutionary change, given the
361 decreased likelihood that they will give rise to harmful pleiotropic effects, especially during
362 development [2, 63].

363 We have previously characterised interspecies differences in human and chimpanzee
364 iPSCs across a suite of regulatory mechanisms – DNA methylation, H3K27ac, H3K27me3 [44],
365 transposable element activity [64] and now, chromatin accessibility and TF binding activity.
366 Though it is difficult to compare divergence across different mechanisms given the multitude of
367 approaches used to collect the data, it seems reasonably safe to state that divergence in all of the
368 regulatory mechanisms we studied is consistently of lesser magnitude than the inter-species
369 changes we observed in gene expression levels. It is unlikely that this is due to consistently
370 worse resolution across this suite of assays than in RNA-seq. Even at the level of RNA-seq, we
371 observed very little intra-specific variation in iPSCs relative to that seen in other somatic tissues,
372 which lead to a dramatic increase in statistical power to identify differential gene expression
373 between the two species [44]. Thus it might be the case that the pluripotent state is subject to
374 much stronger evolutionary constraint than terminally differentiated downstream cell types.

375

376

377

378 **Methods:**

379 *Data collection and sequencing*

380 We generated ATAC-seq libraries from 6 previously described chimpanzee iPSC lines
381 and 6 previously described human iPSC lines [44]. All lines were cultured under feeder free
382 conditions on hESC-grade Matrigel (Corning) and Essential 8 (Gibco) media as previously
383 described. Paired end ATAC-seq libraries were generated as in [43] with a single exception: we
384 collected 200,000 cells per pellet rather than the 50,000 described in the original protocol, solely
385 because this made the pellet visible to the naked eye. Additionally, we collected two pellets at
386 the same time for each cell line in case library preparations failed. Sample collection and library
387 preparation were randomised with respect to species at all times. All libraries were multiplexed
388 together and sequenced to an average depth of roughly 252 million reads across 3 flow cells of
389 an Illumina HiSeq 2500; more details of the sequencing output are available in supplementary
390 table 1.

391

392 *Read mapping, QC, and definition of high-confidence orthologous regions*

393 To compare patterns of chromatin accessibility and TF binding between species we
394 subjected all reads to a series of stringent QC and filtering steps. First, we mapped all reads
395 independently to either the human (hg19) or chimpanzee (panTro 2.1.3) genomes using BWA
396 0.7.9a-r786 [65], allowing a maximum of two mismatches per read and a maximum fragment
397 size of 5000 base pairs for paired-end mates. Reads with mapping quality < 30, unmapped reads,
398 multi-mapping reads and reads that mapped outside the autosomes and X chromosome were
399 discarded. We also discarded all reads that mapped to mitochondrial DNA for the main analyses
400 reported in the text (but see below). Next, we identified and removed reads produced by PCR
401 duplicates with Picard Tools (version 1.129; <http://broadinstitute.github.io/picard/>). The number
402 of reads retained at each step is summarised in supplementary table 1.

403 Additionally, to control for differences in genome assembly quality, and to ensure fair
404 comparisons between the two species, we retained only those reads that fell within a defined set
405 of windows with high orthology between human and chimpanzee. The list of windows was
406 generated by first retrieving the hg19toPanTro3 liftOver chain file from the UCSC Genome
407 Browser [66], and using the regions that are part of the chain to generate non-overlapping 100 bp
408 windows, for a total of 27.2 million 100bp windows. We then used liftOver [67] to test whether
409 each of these windows could be lifted over to a single site in the chimpanzee genome, and from

410 there back to its original location in the human genome, allowing for a maximum 20% change in
411 size during each lift over process. Windows that failed either of these steps were discarded. In
412 parallel, we calculated the uniqueness of every 50-mer in the chimpanzee and human genomes
413 using the GEM suite [68]. Windows where greater than 20% of the 50-mers in either species
414 were not uniquely mappable in their respective genomes, were removed from further analysis. In
415 total, we retained ~17.6 million 100-bp high-quality windows for downstream analyses.

416 Each flow cell was processed separately at this stage, to allow us to identify any flow-cell
417 or library-specific effects. Some libraries yielded very low numbers of mapped reads (< 100,000)
418 and were clear outliers in QC. Although this is probably due to poor multiplexing rather than
419 poor complexity, we discarded these libraries from all downstream analyses and replaced them
420 with new libraries generated from the second pellet collected from the relevant individual. On
421 the whole, we found overall robust reproducibility between libraries prepared from the same
422 individual and thus present analyses at the individual level throughout. Ultimately, between 3.99
423 and 20.8% of generated raw reads per sample met all QC thresholds and were used in all of the
424 following analyses, unless explicitly noted otherwise.

425

426 *Mitochondrial lineage reconstruction*

427 During mapping we observed that a large fraction of reads (between 13.7% and 57.1%)
428 were of mitochondrial origin. Because our chimpanzee iPSC lines are derived from second or
429 third generation captive born individuals from the Yerkes Primate Research Centre, their genetic
430 background is both admixed and unascertained. Therefore, we used MIRA 4 [69] and MITObim
431 1.8 [70] to assemble the mitochondrial genomes of all sequenced libraries. We made use of the
432 publicly available pipeline at <https://github.com/chrishah/MITObim>, using the 'quick' option, and
433 baited the assemblies with the human rCRS ([Genbank: NC_012920](#)) or the chimpanzee reference
434 mtDNA genome ([Genbank: NC_001643](#)), as appropriate. We successfully generated single-
435 contig mtDNA assemblies for all libraries. A maximum likelihood tree of the sequences, built
436 with MEGA 7 [71], confirmed that in all cases except for one chimpanzee library, multiple
437 libraries from the same individual grouped together as sister taxa (supplementary figure 5),
438 confirming there were no sample swaps during sample preparation or sequencing. The one
439 mismatched sample fell into a clade with a maternal cousin. In order to help establish the
440 maternal lineage of the lines in our chimpanzee panel, we also included additional chimpanzee

441 mtDNA sequences from [72] representing all four chimpanzee sub-species (supplementary figure
442 1).

443

444 *Identification of fragmentation biases*

445 As crude indicators of quality of the ATAC-seq data, we examined the fragment size
446 distribution of each library. Although libraries were prepared randomly with regards to species,
447 there is a clear and consistent difference in fragmentation patterns between the two species, with
448 chimpanzee libraries exhibiting an excess of short fragments (<100bp) originating in the
449 nucleosome-free region relative to humans. It is unclear whether this reflects a biological or
450 technical confounder, but to better quantify it we defined a series of *ad hoc* metrics and tallied
451 the number of read pairs with fragment sizes between 50-59 bp, 100-109 bp, 150-159 bp and
452 190-199 bp, as well as the ratios between these measurements, and tested for associations
453 between these measurements and all principal components in the data. In the orthologous TSS
454 data, PC1 is strongly associated with the ratio of fragments 50-59 bp long to those 100-109 bp
455 long (henceforth ratio 50/100, $P = 1.7 \times 10^{-4}$, all values are available in supplementary table 2) and
456 50-59 bp to 150-159 bp (henceforth ratio 50/150, $P = 3.4 \times 10^{-5}$), and PC2 is additionally
457 associated with the ratio of reads mapping to orthologous TSS vs reads mapping to orthologous
458 background regions ($P = 3.7 \times 10^{-4}$), which may be an indicator of data noisiness. No PCs are
459 associated with overall sequencing depth. When we consider the 5,675 'highly accessible' regions
460 the association with fragmentation bias, although diminished, remains significant (ratio 50/100
461 PC1 $P = 4.2 \times 10^{-4}$, ratio 50/150 PC1 $P = 0.04$).

462 We note that cell lines were cultured together by the same individual (IGR), and
463 randomised at all times relative to species to prevent the introduction of technical biases
464 associated with our variable of interest. It is possible that this difference in fragmentation
465 patterns is driven by increased sensitivity of chimpanzee iPSCs to the lysis buffers used in
466 ATAC-seq, or by slightly different responses to our restrictive cell culture conditions. We have
467 systematically elected to be conservative when interpreting our results, although in light of our
468 thorough quality control pipeline we are confident that the trends we observe are reflective of
469 true biological signal. Regardless of cause, this excess of short fragments in chimpanzees leads
470 to a decrease in our ability to differentiate between open and closed chromatin in chimpanzees
471 relative to humans, especially at small scales, and consequently to a loss of power to detect

472 binding in chimpanzees relative to humans. What appears to be a human-unique event is likely to
473 be actually conserved between the two species. For the same reason, those chimpanzee-unique
474 accessibility events we do observe are likely to be true positives.

475 We considered subsampling the chimpanzee data to match the average human fragment
476 distribution, but given both our uncertainty as to the cause of this difference and the relatively
477 low number of reads that are retained given our exhaustive QC strategy, we reasoned that it was
478 more likely to lead to an overall loss, rather than a gain, in power to detect differences.

479

480 *Characterising activity near transcription start sites*

481 To examine patterns of activity near transcription start sites, we began by defining a set
482 of highly orthologous meta-exons between human and chimpanzee as in [21], using Ensemblv75
483 (February 2014; code and documentation for meta-exon identification is available at
484 [http://www.bitbucket.org/ee_reh_neh/orthoexon]; [73]). This data set contains at least 1
485 orthologous meta-exon associated with 40,075 human genes. We then defined the 5'-most
486 position of the first meta-exon of each gene (adjusted for strand direction) as its orthologous
487 transcription start site (orthologous TSS), and computed the number of high-quality orthologous
488 100bp windows (defined above) that fell within a 2kb window centred on the orthologous TSS.
489 We retained only those regions with > 50% overlap, which yielded a set of 28,238 orthologous
490 TSSs. As an additional filtering step, we discarded 7,493 orthologous TSSs that did not span the
491 annotated human TSS for the relevant gene according to Ensembl 75 GENCODE Basic [74]
492 annotations (when multiple GENCODE Basic transcripts were associated with a single gene, we
493 used the average location of all TSS). In total, we retained 20,745 orthologous TSSs.

494 In parallel, we defined a set of randomly selected 2kb orthologous background regions
495 matched to the orthologous TSSs for broad mappability, additionally requiring that they be at
496 least 5kb away from any annotated orthologous TSS. To identify regions with read depths
497 suggestive of activity above the background cutting rate, we used k-means clustering to group
498 the orthologous TSS data into three clusters on the basis of mean CPM by species. We set $k = 3$
499 rather than 2 to capture the long left tail of the distribution, which is readily visible in figure 1a.

500

501 *Differential expression and accessibility analyses*

502 For consistency with our set of defined orthologous TSS we reanalysed previously
503 published RNA-sequencing data from these cell lines using our updated set of orthologous
504 exons. Analyses were performed using the same analysis script as previously [44]; overall we
505 find 4,244 differentially expressed genes out of 12,674 testable genes, and high concordance
506 between the old and new results ($\rho = 0.89$; past results included an additional human and
507 chimpanzee individual). We intersected the expression data with our set of high-quality
508 orthologous TSSs, and found a total of 5,675 genes with evidence for having an accessible
509 orthologous TSS in at least one species - although the majority, 4,644, were accessible in both
510 species. Of those regions seen only in one species, 873 are humans-unique, and 158 chimpanzee-
511 unique. We tested the highly accessible regions for significant interspecies differences with
512 limma [75] and the same steps that we used to test for differential expression. Given the partially
513 confounded covariate structure, we considered four different models:

- 514 1. $y \sim \beta_{\text{species}} + \epsilon$;
- 515 2. $y \sim \beta_{\text{species}} + \beta_{50/100} + \epsilon$;
- 516 3. $y \sim \beta_{\text{species}} + \beta_{50/150} + \epsilon$;
- 517 4. $y \sim \beta_{\text{species}} + \beta_{50/100} + \beta_{50/150} + \epsilon$.

518 To identify the best model, we computed the AIC for each gene under each model, and found
519 that models 1 and 3 were associated with the lowest AIC in roughly the same number of genes
520 (1765 for model 1, 1909 for model 3), suggesting an overall greater suitability to the data. Given
521 this result, we performed differential accessibility testing using model 3. However, we note that
522 all four sets of results are qualitatively similar (supplementary figure 6).

523

524 *Estimating transcription factor binding with msCentipede*

525 In order to determine a set of suitable PWM genomics matches to consider for analysis,
526 we downloaded all 640 PWMs in version 10 of the HOCOMOCO database [54]. To ensure fair
527 comparisons between human and chimpanzee data, we scanned both genomes for matches to
528 each PWM, retaining any site with a PWM score ≥ 7 . The adoption of this permissively low
529 threshold was motivated by our desire to capture turnover at the PWM site between the two
530 species. We then used liftOver [67] to ensure that identified human PWM matches were also
531 present in the chimpanzee genome, and could be lifted over back to their original location in the
532 human genome. We then used the k-mer mappability data described above to ensure that at least

533 80% of the 50-mers originating \pm 100bp around the PWM site were uniquely mapping in both
534 species. Finally, we excluded all PWMs associated with transcriptions factors that were either
535 not present or not expressed in at least half of the individuals from one species according to our
536 RNA-seq data, with the exception of master pluripotency regulators OCT4 and NANOG, where
537 we confirmed expression through qPCR. These two genes cannot be reliably assayed through
538 RNA-seq due to the existence of closely related pseudogenes that confound mapping. We also
539 included the *REXI* (also known as *ZFP42*) PWM from HOCOMOCO v11, as we had previously
540 identified it as the sole master regulator of pluripotency that was differentially expressed
541 between the two species. After all of these filtering steps, we retained 133,103,977 possible
542 PWM sites across 306 TFs, 906,612 of which had a PWM score \geq 12 in at least one species.

543 Finally, we used msCentipede [45], an extension of the CENTIPEDE algorithm [46] that
544 can better capture heterogeneity around binding sites, to predict TF binding in all PWM sites in
545 human and chimpanzee iPSCs. The vast majority of log posterior odds distributions reported by
546 msCentipede are unimodal, with an extremely long rightward tail and without a clear break point
547 (supplementary figure 7); chimpanzee distributions in particular are clearly narrower and left-
548 shifted relative to human distributions, again reflecting our diminished power in chimpanzees.
549 As such, we conservatively chose to call a site bound in a species if the log posterior odds were \geq
550 $\log(0.99/0.01)$ at that particular site, and unbound if they fell below that threshold. However, we
551 note that our results are qualitatively robust to lowering this threshold to $\log(0.95/0.05)$ in
552 chimpanzees to account for the decreased power in that species, resulting in an increase in the
553 mean fraction of bound sites in chimpanzee, to 12.1%, a slight decrease in the fraction of sites
554 called as bound in chimpanzee that are also bound in human, to 77.1%, and overall corroboration
555 of our findings.

556 Additionally, HOCOMOCO summarises the quality and predictive value of PWMs using
557 a simple A-B-C-D quality score, which is assigned to each PWM on the basis of receiver
558 operating curve analyses [54]. Nearly half of the PWMs in the dataset ($n = 133$) have been
559 assigned to class D, suggesting they might be information-poor and only capture a small subset
560 of binding activity. This reflects broader observations that TFs can vary dramatically in their
561 degree of preference for specific binding motifs [76]. We do find that the fraction of MPBS
562 predicted to be bound does vary across PWM qualities, but only when we consider binding
563 predictions in the 'high-quality' motif subset in chimpanzees (ANOVA P value = 0.016 genome-

564 wide, 0.03 within 5kb of annotated orthologous TSS). In both of these cases, the significance is
565 driven by differences between class A and D PWMs, but its overall impact appears relatively
566 minor (supplementary figure 8). In light of these findings, we retained all PWMs for downstream
567 analyses.

568 All analyses were performed using R 3.2.2 [77].

569

570 *Data accessibility*

571 All raw ATAC-seq reads have been deposited in GEO/SRA under series number
572 GSE122319, alongside individual msCentipede calls for all high quality sites and summary
573 statistics for all tested PWMs. Previously published RNA-seq data from human and chimpanzee
574 iPSCs is available under BioProject PRJNA260053. For comparisons with published ENCODE
575 data we used all publicly released "optimal idr thresholded peaks" associated with either the
576 human embryonic stem cell line H1 or the human induced pluripotent stem cell line GM23339
577 available on the ENCODE browser (<https://www.encodeproject.org>) as of the 15th of May, 2018.
578 A full list of accession identifiers is provided in supplementary table 6.

579

580

581 **Acknowledgements:**

582 We thank members of the Gilad lab for helpful discussions, as well as Anil Raj and Heejung
583 Shim for advice regarding msCentipede. IGR was supported in part by a Sir Henry Wellcome
584 Postdoctoral Fellowship. SG was supported by Marie Skłodowska-Curie actions grant 655732-
585 WhereWolf. This work was additionally supported by an NIGMS award to YG.

586

587 **References**

- 588
- 589 1. King, M. and A. Wilson, *Evolution at two levels in humans and chimpanzees*. Science,
590 1975. **188**(4184): p. 107-116.
 - 591 2. Carroll, S.B., *Evolution at two levels: On genes and form*. PLoS Biology, 2005. **3**(7): p.
592 1159-1166.
 - 593 3. Britten, R.J. and E.H. Davidson, *Repetitive and non-repetitive DNA sequences and a*
594 *speculation on the origins of evolutionary novelty*. Quarterly Review of Biology, 1971.
595 **46**(2): p. 111-138.
 - 596 4. Jacob, F., *Evolution and tinkering*. Science, 1977. **196**(4295): p. 1161-6.
 - 597 5. Gallego Romero, I., I. Ruvinsky, and Y. Gilad, *Comparative studies of gene expression*
598 *and the evolution of gene regulation*. Nature Reviews Genetics, 2012. **13**(7): p. 505-516.
 - 599 6. Bejerano, G., et al., *Ultraconserved elements in the human genome*. Science, 2004.
600 **304**(5675): p. 1321-5.
 - 601 7. Siepel, A., et al., *Evolutionarily conserved elements in vertebrate, insect, worm, and*
602 *yeast genomes*. Genome Research, 2005. **15**(8): p. 1034-50.
 - 603 8. Prabhakar, S., et al., *Accelerated evolution of conserved noncoding sequences in humans*.
604 Science, 2006. **314**(5800): p. 786.
 - 605 9. Pollard, K.S., et al., *Forces shaping the fastest evolving regions in the human genome*.
606 PLoS Genetics, 2006. **2**(10): p. e168.
 - 607 10. Booker, B.M., et al., *Bat Accelerated Regions Identify a Bat Forelimb Specific Enhancer*
608 *in the HoxD Locus*. PLoS Genetics, 2016. **12**(3): p. e1005738-e1005738.
 - 609 11. Eckalbar, W.L., et al., *Transcriptomic and epigenomic characterization of the developing*
610 *bat wing*. Nature Genetics, 2016. **48**(5): p. 528-536.
 - 611 12. Pollard, K., et al., *An RNA gene expressed during cortical development evolved rapidly in*
612 *humans*. Nature, 2006. **443**(7108): p. 167-172.
 - 613 13. Hubisz, M.J. and K.S. Pollard, *Exploring the genesis and functions of Human*
614 *Accelerated Regions sheds light on their role in human evolution*. Current Opinion in
615 Genetics and Development, 2014. **29**: p. 15-21.
 - 616 14. Kvon, E.Z., et al., *Progressive Loss of Function in a Limb Enhancer during Snake*
617 *Evolution*. Cell, 2016. **167**(3): p. 633-642.e11.

- 618 15. O'Brown, N.M., et al., *A recurrent regulatory change underlying altered expression and*
619 *Wnt response of the stickleback armor plates gene EDA*. *Elife*, 2015. **4**: p. e05290.
- 620 16. Reilly, S., et al., *Evolutionary changes in promoter and enhancer activity during human*
621 *corticogenesis*. *Science*, 2015. **347**(6226): p. 1155-1159.
- 622 17. Bakken, T.E., et al., *A comprehensive transcriptional map of primate brain development*.
623 *Nature*, 2016. **535**(7612): p. 367-375.
- 624 18. Cotney, J., et al., *The Evolution of Lineage-Specific Regulatory Activities in the Human*
625 *Embryonic Limb*. *Cell*, 2013. **154**(1): p. 185-196.
- 626 19. Prabhakar, S., et al., *Human-Specific Gain of Function in a Developmental Enhancer*.
627 *Science*, 2008. **321**(5894): p. 1346-1350.
- 628 20. Blekhman, R., et al., *Gene regulation in primates evolves under tissue-specific selection*
629 *pressures*. *PLoS Genetics*, 2008. **4**(11): p. e1000271-e1000271.
- 630 21. Blekhman, R., et al., *Sex-specific and lineage-specific alternative splicing in primates*.
631 *Genome Research*, 2010. **20**(2): p. 180-189.
- 632 22. Brawand, D., et al., *The evolution of gene expression levels in mammalian organs*.
633 *Nature*, 2011. **478**(7369): p. 343-348.
- 634 23. Khaitovich, P., et al., *Regional patterns of gene expression in human and chimpanzee*
635 *brains*. *Genome Research*, 2004. **14**(8): p. 1462-1473.
- 636 24. Enard, W., et al., *Intra- and interspecific variation in primate gene expression patterns*.
637 *Science*, 2002. **296**(5566): p. 340-3.
- 638 25. He, Z., et al., *Comprehensive transcriptome analysis of neocortical layers in humans,*
639 *chimpanzees and macaques*. *Nature Neuroscience*, 2017.
- 640 26. Cain, C.E., et al., *Gene expression differences among primates are associated with*
641 *changes in a histone epigenetic modification*. *Genetics*, 2011. **187**(4): p. 1225-1234.
- 642 27. Zhou, X., et al., *Epigenetic modifications are associated with inter-species gene*
643 *expression variation in primates*. *Genome Biology*, 2014. **15**(12): p. 547-547.
- 644 28. Pai, A., et al., *A Genome-Wide Study of DNA Methylation Patterns and Gene Expression*
645 *Levels in Multiple Human and Chimpanzee Tissues*. *PLOS Genetics*, 2011. **7**(2): p.
646 e1001316-e1001316.
- 647 29. Hernando-Herraez, I., et al., *The interplay between DNA methylation and sequence*
648 *divergence in recent human evolution*. *bioRxiv*, 2015: p. 015966-015966.

- 649 30. Hernando-Herraez, I., et al., *Dynamics of DNA Methylation in Recent Human and Great*
650 *Ape Evolution*. PLoS Genetics, 2013. **9**(9): p. e1003763-e1003763.
- 651 31. Bilgin Sonay, T., et al., *Tandem repeat variation in human and great ape populations*
652 *and its impact on gene expression divergence*. Genome Research, 2015. **25**(11): p. 1591-
653 9.
- 654 32. Blekhman, R., A. Oshlack, and Y. Gilad, *Segmental duplications contribute to gene*
655 *expression differences between humans and chimpanzees*. Genetics, 2009. **182**(2): p. 627-
656 630.
- 657 33. Ward, M.C., et al., *Silencing of transposable elements may not be a major driver of*
658 *regulatory evolution in primate iPSCs*. Elife, 2018. **7**.
- 659 34. Nowick, K., et al., *Differences in human and chimpanzee gene expression patterns define*
660 *an evolving network of transcription factors in brain*. Proc Natl Acad Sci U S A, 2009.
661 **106**(52): p. 22358-63.
- 662 35. Shibata, Y., et al., *Extensive Evolutionary Changes in Regulatory Element Activity during*
663 *Human Origins Are Associated with Altered Gene Expression and Positive Selection*.
664 PLoS Genetics, 2012. **8**(6): p. e1002789-e1002789.
- 665 36. Berthelot, C., et al., *Complexity and conservation of regulatory landscapes underlie*
666 *evolutionary resilience of mammalian gene expression*. Nat Ecol Evol, 2018. **2**(1): p.
667 152-163.
- 668 37. Paris, M., et al., *Extensive divergence of transcription factor binding in Drosophila*
669 *embryos with highly conserved gene expression*. PLoS Genet, 2013. **9**(9): p. e1003748.
- 670 38. Schmidt, D., et al., *Five-vertebrate ChIP-seq reveals the evolutionary dynamics of*
671 *transcription factor binding*. Science, 2010. **328**(5981): p. 1036-40.
- 672 39. Wilson, M.D., et al., *Species-specific transcription in mice carrying human chromosome*
673 *21*. Science, 2008. **322**(5900): p. 434-8.
- 674 40. Li, X.Y., et al., *The role of chromatin accessibility in directing the widespread,*
675 *overlapping patterns of Drosophila transcription factor binding*. Genome Biol, 2011.
676 **12**(4): p. R34.
- 677 41. Peng, P.-C., et al., *Evolutionary changes in DNA accessibility and sequence predict*
678 *divergence of transcription factor binding and enhancer activity*. bioRxiv, 2018.

- 679 42. Vierstra, J. and J.A. Stamatoyannopoulos, *Genomic footprinting*. Nature Methods, 2016.
680 **13**(3): p. 213-221.
- 681 43. Buenrostro, J.D., et al., *Transposition of native chromatin for fast and sensitive*
682 *epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position*.
683 Nature Methods, 2013. **10**(12): p. 1213-1218.
- 684 44. Gallego Romero, I., et al., *A panel of induced pluripotent stem cells from chimpanzees: a*
685 *resource for comparative functional genomics*. eLife, 2015. **4**: p. e07103-e07103.
- 686 45. Raj, A., et al., *msCentipede: Modeling Heterogeneity across Genomic Sites and*
687 *Replicates Improves Accuracy in the Inference of Transcription Factor Binding*. PLOS
688 ONE, 2015. **10**(9): p. e0138030-e0138030.
- 689 46. Pique-Regi, R., et al., *Accurate inference of transcription factor binding from DNA*
690 *sequence and chromatin accessibility data*. Genome Research, 2011. **21**(3): p. 447-455.
- 691 47. Ashburner, M., et al., *Gene Ontology: tool for the unification of biology*. Nature Genetics,
692 2000. **25**(1): p. 25-29.
- 693 48. Rada-Iglesias, A., et al., *A unique chromatin signature uncovers early developmental*
694 *enhancers in humans*. Nature, 2010. **470**(7333): p. 279-283.
- 695 49. Roadmap Epigenomics, C., et al., *Integrative analysis of 111 reference human*
696 *epigenomes*. Nature, 2015. **518**(7539): p. 317-30.
- 697 50. Bu, L., et al., *Human ISL1 heart progenitors generate diverse multipotent cardiovascular*
698 *cell lineages*. Nature, 2009. **460**(7251): p. 113-7.
- 699 51. McCormick, M.B., et al., *NeuroD2 and neuroD3: distinct expression patterns and*
700 *transcriptional activation potentials within the neuroD gene family*. Mol Cell Biol, 1996.
701 **16**(10): p. 5792-800.
- 702 52. Masserdotti, G., et al., *Transcriptional Mechanisms of Proneural Factors and REST in*
703 *Regulating Neuronal Reprogramming of Astrocytes*. Cell Stem Cell, 2015. **17**(1): p. 74-
704 88.
- 705 53. Tsuda, H., et al., *Structure and promoter analysis of Math3 gene, a mouse homolog of*
706 *Drosophila proneural gene atonal. Neural-specific expression by dual promoter*
707 *elements*. J Biol Chem, 1998. **273**(11): p. 6327-33.

- 708 54. Kulakovskiy, I.V., et al., *HOCOMOCO: Expansion and enhancement of the collection of*
709 *transcription factor binding sites models*. Nucleic Acids Research, 2016. **44**(D1): p.
710 D116-D125.
- 711 55. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human genome*.
712 Nature, 2012. **489**(7414): p. 57-74.
- 713 56. Young, R.A., *Control of the Embryonic Stem Cell State*. Cell, 2011. **144**(6): p. 940-954.
- 714 57. Marchetto, M.C.N., et al., *Differential L1 regulation in pluripotent stem cells of humans*
715 *and apes*. Nature, 2013. **503**(7477): p. 525-529.
- 716 58. Chronis, C., et al., *Cooperative Binding of Transcription Factors Orchestrates*
717 *Reprogramming*. Cell, 2017. **168**(3): p. 442-459 e20.
- 718 59. Orkin, S.H. and K. Hochedlinger, *Chromatin Connections to Pluripotency and Cellular*
719 *Reprogramming*. Cell, 2011. **145**(6): p. 835-850.
- 720 60. Ng, H.-H. and M. Surani, *The transcriptional and signalling networks of pluripotency*.
721 Nat Cell Biol, 2011. **13**(5): p. 490-496.
- 722 61. Weirauch, M.T., et al., *Determination and inference of eukaryotic transcription factor*
723 *sequence specificity*. Cell, 2014. **158**(6): p. 1431-1443.
- 724 62. Bradley, R.K., et al., *Binding site turnover produces pervasive quantitative changes in*
725 *transcription factor binding between closely related Drosophila species*. PLoS Biol,
726 2010. **8**(3): p. e1000343.
- 727 63. Wray, G., *The evolutionary significance of cis-regulatory mutations*. Nature Reviews
728 Genetics, 2007. **8**(3): p. 206-216.
- 729 64. Ward, M.C., et al., *Silencing of transposable elements may not be a major driver of*
730 *regulatory evolution in primate induced pluripotent stem cells*. bioRxiv, 2017: p. 1-46.
- 731 65. Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows-Wheeler*
732 *transform*. Bioinformatics, 2009. **25**(14): p. 1754-1760.
- 733 66. Kent, W.J., et al., *The human genome browser at UCSC*. Genome Res, 2002. **12**(6): p.
734 996-1006.
- 735 67. Hinrichs, A.S., et al., *The UCSC Genome Browser Database: update 2006*. Nucleic
736 Acids Res, 2006. **34**(Database issue): p. D590-8.
- 737 68. Derrien, T., et al., *Fast computation and applications of genome mappability*. PLoS
738 ONE, 2012. **7**(1): p. e30377-e30377.

- 739 69. Chevreux, B., T. Wetter, and S. Suhai, *Genome Sequence Assembly Using Trace Signals*
740 *and Additional Sequence Information*, in *Computer Science and Biology: Proceedings of*
741 *the German Conference on Bioinformatics*. 1999. p. 45-56.
- 742 70. Hahn, C., L. Bachmann, and B. Chevreux, *Reconstructing mitochondrial genomes*
743 *directly from genomic next-generation sequencing reads - A baiting and iterative*
744 *mapping approach*. *Nucleic Acids Research*, 2013. **41**(13): p. e129--e129.
- 745 71. Tamura, K., et al., *MEGA5: molecular evolutionary genetics analysis using maximum*
746 *likelihood, evolutionary distance, and maximum parsimony methods*. *Mol Biol Evol*,
747 2011. **28**(10): p. 2731-9.
- 748 72. Bjork, A., et al., *Evolutionary history of chimpanzees inferred from complete*
749 *mitochondrial genomes*. *Molecular Biology and Evolution*, 2011. **28**(1): p. 615-623.
- 750 73. Yates, A., et al., *Ensembl 2016*. *Nucleic Acids Res*, 2016. **44**(D1): p. D710-6.
- 751 74. Harrow, J., et al., *GENCODE: the reference human genome annotation for The*
752 *ENCODE Project*. *Genome Res*, 2012. **22**(9): p. 1760-74.
- 753 75. Ritchie, M., et al., *limma powers differential expression analyses for RNA-sequencing*
754 *and microarray studies*. *Nucleic Acids Research*, 2015. **43**(7): p. e47--e47.
- 755 76. Karimzadeh, M. and M.M. Hoffman, *Virtual ChIP-seq: predicting transcription factor*
756 *binding by learning from the transcriptome*. bioRxiv, 2018.
- 757 77. R Core Team, *R: A Language and Environment for Statistical*
758 *Computing*, R Foundation for Statistical Computing, Editor. 2015: Vienna.
- 759
- 760
- 761
- 762
- 763

764 **Tables:**

765 Supplementary table 1: Sequencing depths and quality filtering steps.

766

767 Supplementary table 2: P-values for the association between each principal component and
768 possible covariates

769

770 Supplementary table 3: Results of differential accessibility testing across 5,675 'highly
771 accessible' orthologous TSS using model 3 as described above.

772

773 Supplementary table 4: Results of differential expression testing across 12,674 orthologous
774 genes.

775

776 Supplementary table 5: List of genes driving the GO enrichment results.

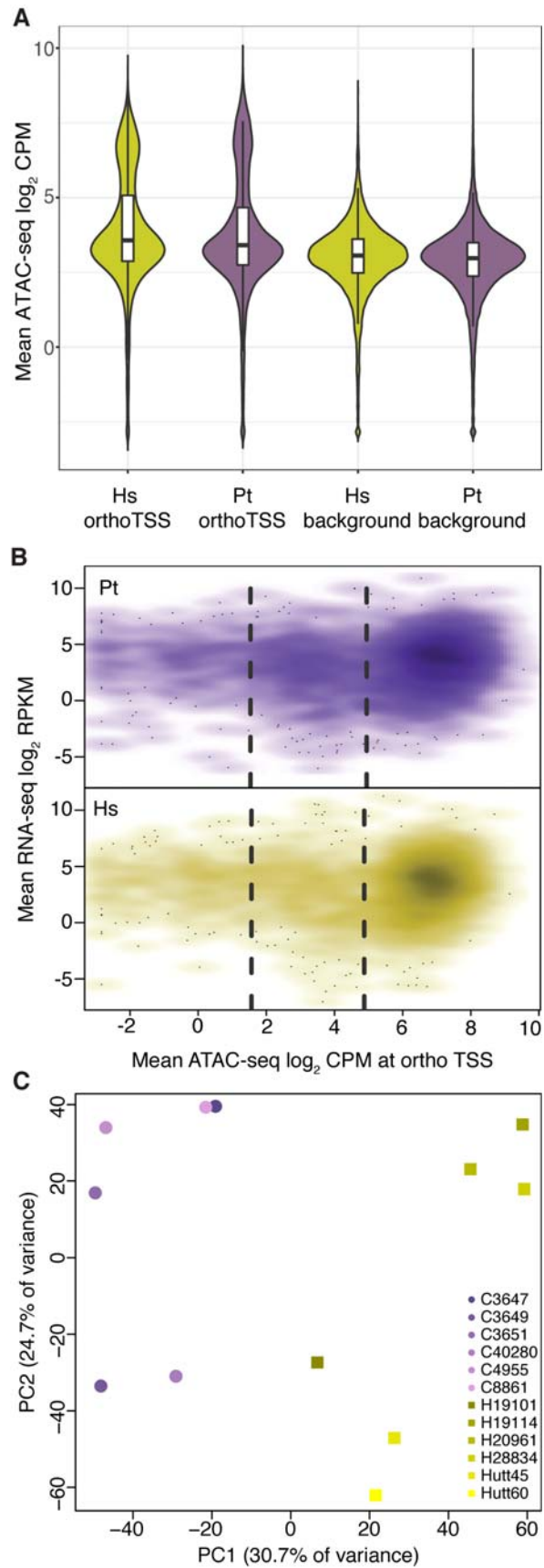
777

778 Supplementary table 6: Summary of ENCODE TF ChIP datasets used in this publication.

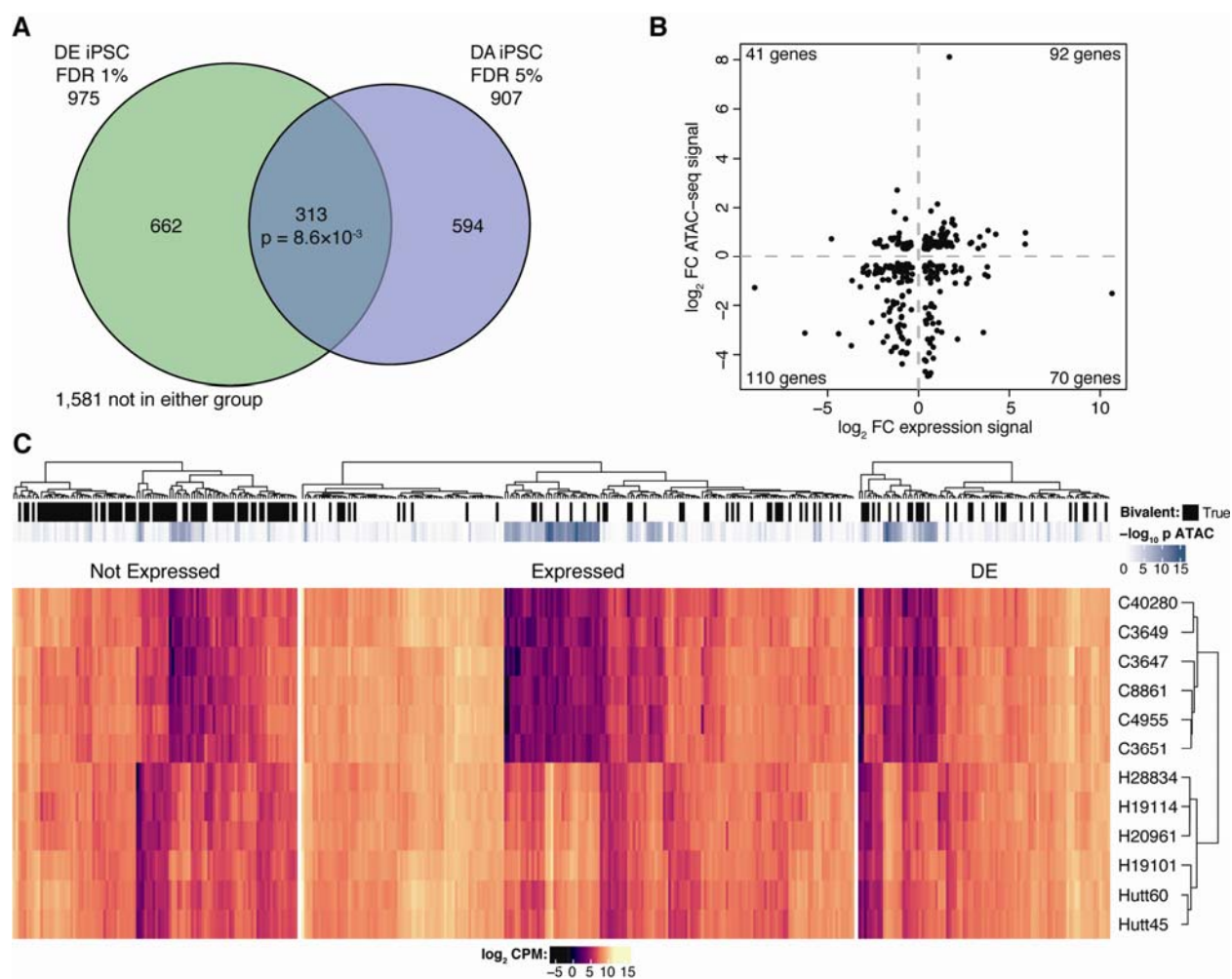
779

780

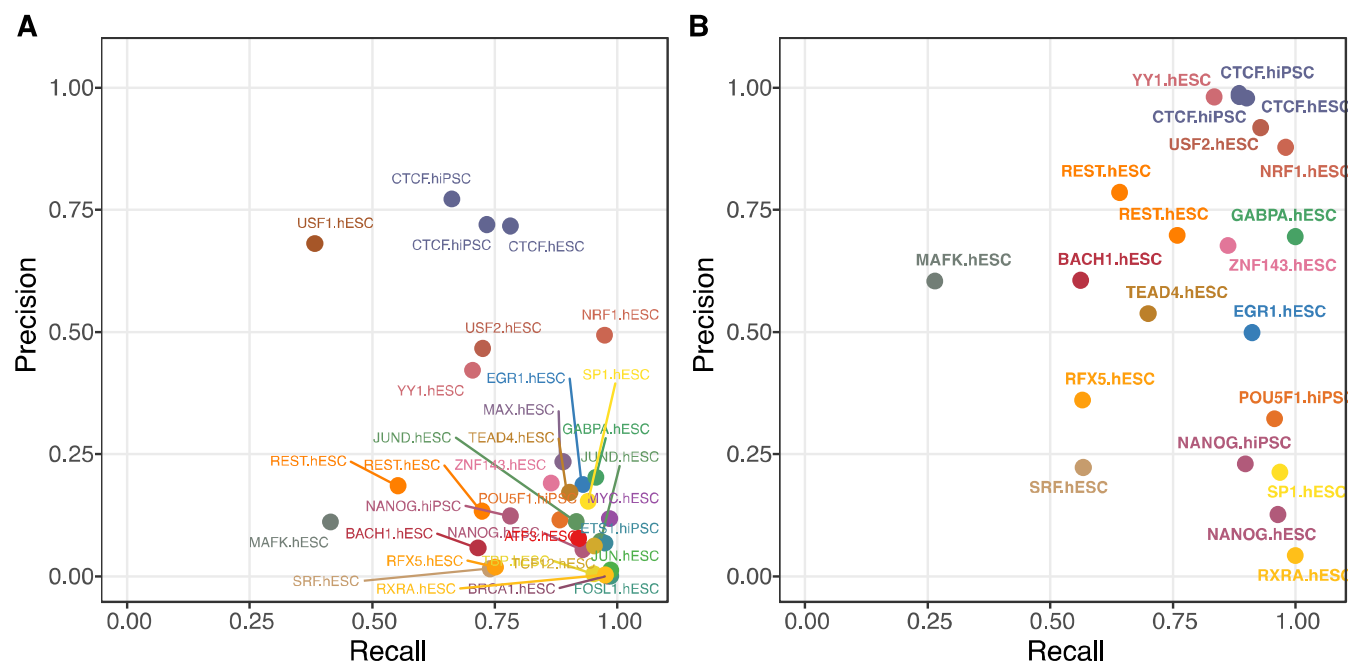
781



783 **Figure 1: Patterns of chromatin accessibility in human and chimpanzee iPSCs. a.**
784 distribution of mean log₂ CPM across 20,745 orthologous transcription start sites and 20,745
785 orthologous background regions. **b.** Expression (log₂ RPKM) and accessibility (log₂ CPM) in
786 humans around the orthologous transcription start site of 4,210 genes in the highly accessible
787 orthologous TSS subset expressed in chimpanzee (top) or human (bottom) iPSCs. The dashed
788 lines indicate cluster boundaries identified through k-means clustering (k = 3) as described in the
789 text. **c.** Principal component analysis of 5,675 orthologous transcription start sites.



790
791 **Figure 2: Testing for differentially accessible regions uncovers biologically meaningful**
792 **signal.** **a.** Overlap between differentially expressed and differentially accessible genes. p value
793 from hypergeometric test. **b.** Effect directions across the 313 genes that are both differentially
794 accessible and differentially expressed between human and chimpanzee iPSCs. **c.** Heatmap of
795 394 genes associated with significant GO terms, grouped by their expression status.
796



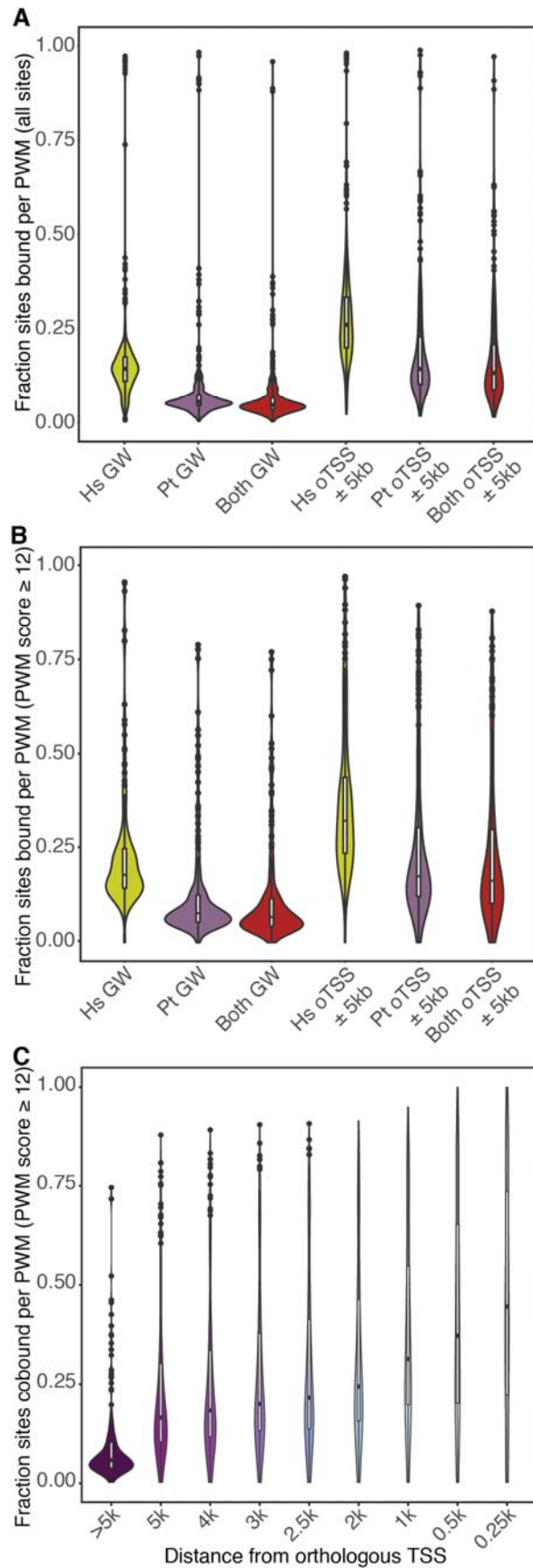
797

798 **Figure 3: Recall and precision for msCentipede human binding predictions compared**

799 **against ENCODE ChIP-seq data. a.** At all sites with a PWM score ≥ 7 . **b.** Within the subset of

800 motifs with a PWM score ≥ 12 .

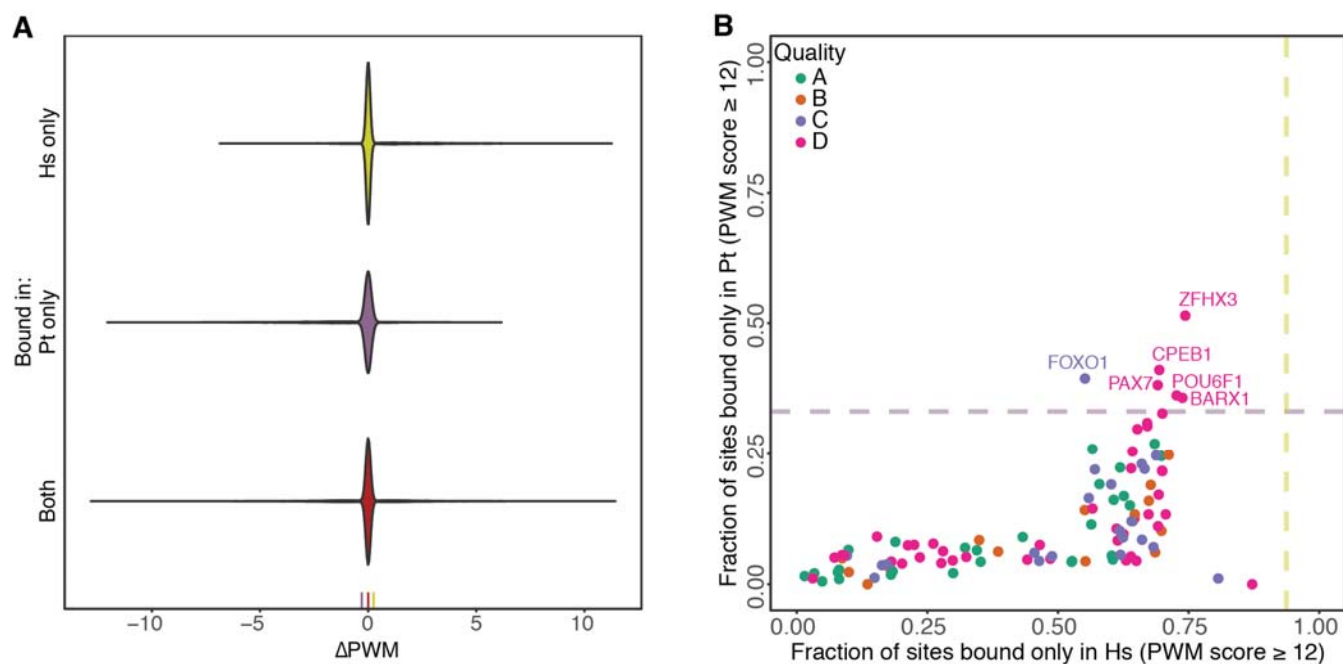
801



803 **Figure 4: Transcription factor binding activity inferred by msCentipede in human and**
804 **chimpanzee iPSCs. a.** At all sites with a PWM score ≥ 7 . **b.** Within the subset of motifs with a
805 PWM score ≥ 12 . **c.** Fraction of sites bound in both species at decreasing distance from
806 orthologous TSS.

807

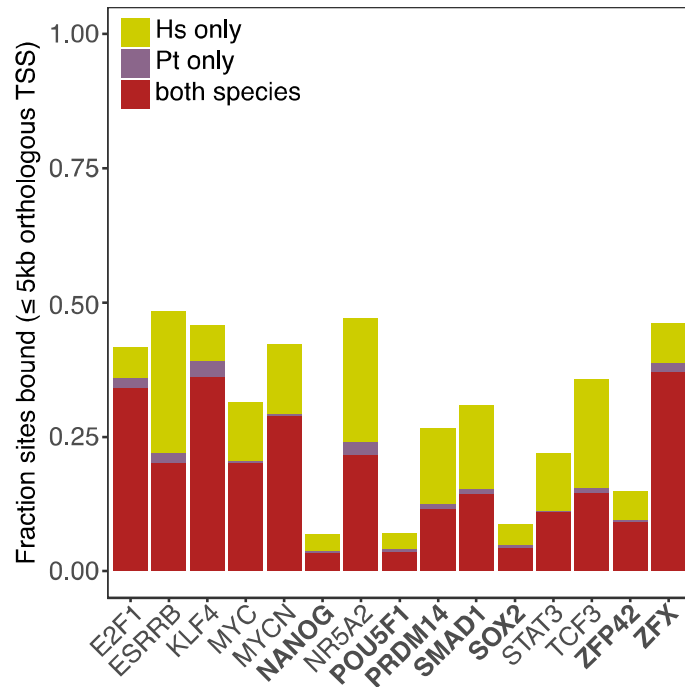
808



809

810 **Figure 5: Binding differences between species a.** Distribution of Δ PWM scores by binding
811 status. The three marks at the bottom highlight the median value for each dataset. **b.** Species-
812 unique binding by PWM. The dashed purple and yellow lines denote 2 standard deviations away
813 from the species-wide average for chimpanzee and humans, respectively.

814



815

816 **Figure 6: Predicted binding activity across 15 key pluripotency regulators genome-wide.**

817 Factors in boldface are associated with at least 50 sites with a PWM score ≥ 12 .

818

819 **Supplementary figures:**

820 Supplementary figure 1: Maximum likelihood tree from reconstructed and publicly available
821 chimpanzee mtDNA sequences.

822

823 Supplementary figure 2: Principal component analysis of a. 20,745 orthologous transcription
824 start sites genome-wide. b. 20,745 control regions genome-wide.

825

826 Supplementary figure 3: Log_2 cpm distributions by individual around orthologous TSS and
827 orthologous background regions.

828

829 Supplementary figure 4: Fraction of sites in the high-quality motif subset predicted to be bound
830 within 5kb of an annotated orthologous TSS and genome wide.

831

832 Supplementary figure 5: Maximum likelihood tree from reconstructed human and chimpanzee
833 mtDNA sequences.

834

835 Supplementary figure 6: overlap in differential accessibility results across the four considered
836 models.

837

838 Supplementary figure 7: Log posterior binding probabilities for all 306 PWMs in the study, by
839 species. The dashed vertical line indicates $\ln(0.99/0.01)$, our threshold for calling a site bound.

840

841 Supplementary figure 8: Fraction of sites predicted to be bound by PWM across the four main
842 HOCOMOCO quality classes. Pairwise test P-values are reported after Tukey's post hoc HSD
843 correction.

844