



ATAC-STARR-seq reveals transcription factor-bound activators and silencers across the chromatin accessible human genome

Tyler J. Hansen and Emily Hodges

Genome Res. published online July 20, 2022

Access the most recent version at doi:[10.1101/gr.276766.122](https://doi.org/10.1101/gr.276766.122)

P<P	Published online July 20, 2022 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **ATAC-STARR-seq reveals transcription factor-bound activators and silencers across**
2 **the chromatin accessible human genome**

3 Tyler J. Hansen¹ and Emily Hodges^{1,2,*}

4 ¹ Department of Biochemistry, Vanderbilt University School of Medicine, Nashville, TN,
5 37232, USA

6 ² Vanderbilt Genetics Institute, Vanderbilt University School of Medicine, Nashville, TN,
7 37232, USA

8 * Correspondence: Tel: +1 615 875 9991; Email: emily.hodges@vanderbilt.edu

9 **Running Title (50-character limit):** Dissecting human gene regulation with ATAC-STARR

10 **ABSTRACT**

11 Massively parallel reporter assays (MPRAs) test the capacity of putative gene regulatory
12 elements to drive transcription on a genome-wide scale. Most gene regulatory activity occurs
13 within accessible chromatin, and recently described methods have combined assays that
14 capture these regions—such as assay for transposase-accessible chromatin using
15 sequencing (ATAC-seq)—with self-transcribing active regulatory region sequencing
16 (STARR-seq) to selectively assay the regulatory potential of accessible DNA (ATAC-
17 STARR-seq). Here, we report an integrated approach that quantifies activating and silencing
18 regulatory activity, chromatin accessibility, and transcription factor (TF) occupancy with one
19 assay using ATAC-STARR-seq. Our strategy, including important updates to the ATAC-
20 STARR-seq assay and workflow, enabled high-resolution testing of ~50 million unique DNA
21 fragments tiling ~101,000 accessible chromatin regions in human lymphoblastoid cells. We
22 discovered that 30% of all accessible regions contain an activator, a silencer or both.
23 Although few MPRA studies have explored silencing activity, we demonstrate silencers
24 occur at similar frequencies to activators, and they represent a distinct functional group
25 enriched for unique TF motifs and repressive histone modifications. We further show that
26 Tn5 cut-site frequencies are retained in the ATAC-STARR plasmid library compared to
27 standard ATAC-seq, enabling TF occupancy to be ascertained from ATAC-STARR data.
28 With this approach, we found that activators and silencers cluster by distinct TF footprint
29 combinations and these groups of activity represent different gene regulatory networks of
30 immune cell function. Altogether, these data highlight the multi-layered capabilities of ATAC-
31 STARR-seq to comprehensively investigate the regulatory landscape of the human genome
32 all from a single DNA fragment source.

33 **INTRODUCTION**

34 Transcription is regulated by transcription factors (TFs) and the DNA sequences they bind,
35 called *cis*-regulatory elements. Enhancers, which are a class of *cis*-regulatory elements, are

36 distally located from the genes they target and serve as key drivers of cell-type specific gene
37 expression (Heinz et al. 2015). Because enhancers require TF binding, they are largely
38 dependent on chromatin accessibility to elicit transcriptional activity. Therefore, chromatin
39 accessibility is a vital regulator of enhancer function, and this is evidenced by the
40 observation that ~94% of all ENCODE TF ChIP-seq peaks fall within accessible chromatin
41 (Klemm et al. 2019). In any given cell type, only a small fraction (~2%) of the genome is
42 accessible to TF binding (Thurman et al. 2012; Klemm et al. 2019). In this way, most
43 enhancers are inaccessible and are less likely to drive transcription endogenously.

44 Enhancers are difficult to identify and validate because they lack uniform features and are
45 less constrained by gene proximity than promoters (Gasparini et al. 2020). Massively parallel
46 reporter assays (MPRAs) were developed to test the regulatory potential of thousands to
47 millions of DNA sequences in parallel, providing high-throughput identification of putative
48 enhancers. Overall, MPRAs test the regulatory potential of genomic regions by cloning them
49 *en masse* into a reporter plasmid and leveraging high-throughput sequencing to quantify
50 regulatory activity (Santiago-Algarra et al. 2017). Among the variety of different vector
51 backbones and assay designs applied to MPRAs, Self-Transcribing Active Regulatory
52 Region sequencing (STARR-seq) is uniquely designed to assay an entire genome for
53 regulatory activity (Melnikov et al. 2012; Patwardhan et al. 2012; Arnold et al. 2013; Inoue et
54 al. 2017; Maricque et al. 2017; Muerdter et al. 2018; Kircher et al. 2019). STARR-seq
55 quantifies regulatory activity genome-wide by cloning randomly fragmented genomic DNA
56 into the 3'UTR of the reporter plasmid. Thus, active enhancers drive transcription of
57 themselves, and activity is quantified by the abundance of its own sequence in the transcript
58 pool, removing the need for barcodes that some MPRAs employ. One major limitation of
59 STARR-seq is that it is technically challenging to accommodate the massive size of the
60 human genome; it requires large-scale cloning procedures and produces shallow
61 sequencing coverage of human regulatory elements (Johnson et al. 2018). In addition,
62 STARR-seq assays both accessible and inaccessible chromatin. Thus, many assayed

63 regions are derived from heterochromatin and are less likely to be transcriptionally active in
64 the cell type in question.

65 To narrow the scope of the assay, recent methods have combined STARR-seq with
66 techniques that capture accessible chromatin to specifically test the regulatory potential of
67 accessible DNA (Buenrostro et al. 2013; Wang et al. 2018; Chaudhri et al. 2020; Glaser et al.
68 2021). As a result, these methods only sample a fraction of the human genome (~2%) while
69 assaying nearly all regulatory elements capable of driving transcription endogenously,
70 because they are derived from open chromatin. This approach remains comprehensive while
71 enabling deeper sequencing coverage of biologically relevant genomic regions. Furthermore,
72 integrated approaches have recently been described that combine measurements of
73 chromatin accessibility with analysis of transcription and other epigenomic features from a
74 single population of cells (Kelly et al. 2012; Clark et al. 2018; Barnett et al. 2020; Chen et al.
75 2022). Similarly, ATAC-STARR-seq has the potential to reveal multiple levels of gene
76 regulatory information simultaneously, but this potential has not been explored. In addition, a
77 complete understanding of gene regulatory activity is lacking with most MPRA approaches
78 because silencing activity is largely overlooked, with a few recent exceptions (Doni Jayavelu
79 et al. 2020; Pang and Snyder 2020; Kim et al. 2021); this is potentially due to technical
80 caveats of distinguishing silencers from either that of missing data or interference from head-
81 on transcriptional conflicts or post-transcriptional silencing mechanisms.

82 Here, we demonstrate a new workflow that substantially expands the capabilities of ATAC-
83 STARR-seq to extract and measure gene regulatory information. Using this approach, we
84 aimed to identify both activators and silencers, as well as to simultaneously profile chromatin
85 accessibility, and perform TF footprinting. From a single ATAC-STARR-seq dataset, a multi-
86 layered, integrated view of the human genome can be captured—a feature that has not been
87 explored previously. We provide a protocol and code repository so that this new ATAC-
88 STARR-seq workflow may be easily used and adopted by the field.

89 RESULTS

90 ATAC-STARR-seq Experimental Design

91 The ATAC-STARR-seq approach is divided into the three main parts: 1) ATAC-STARR-seq
92 plasmid library generation, 2) reporter assay, and 3) data analysis (Figure 1A). To generate
93 ATAC-STARR-seq plasmid libraries, nuclei are isolated from a cell type of interest and
94 exposed to Tn5, the cut-and-paste transposase used in the ATAC-seq method (Buenrostro
95 et al. 2013). Tn5 simultaneously cleaves DNA fragments within accessible chromatin and
96 attaches customizable sequence adapters to their 5' ends. ATAC-STARR-seq adapters are
97 designed to serve as homology arms for direct Gibson cloning into the STARR-seq reporter
98 plasmid, which enables cloning of accessible DNA fragments *en masse*. The resulting
99 ATAC-STARR-seq plasmid library consists of millions of unique plasmids each harbouring
100 their own unique open chromatin-derived DNA fragment.

101 In our updated ATAC-STARR-seq workflow, we employ the STARR-seq Ori backbone,
102 where the origin of replication (Ori) functions as the minimal promoter (Muerdter et al. 2018)
103 (Supplemental Table S1). Each plasmid in the ATAC-STARR-seq plasmid library contains a
104 truncated GFP (trGFP) coding sequence, a poly-adenylation signal sequence, the Ori, and
105 the unique accessible DNA fragment being assayed (Figure 1B). Critically, the accessible
106 region is cloned into the 3' UTR, so if the accessible region is active, it interacts with the Ori
107 to drive self-transcription. Thus, an accessible region's level of activity is reflected by its own
108 level of expression. Transcripts from ATAC-STARR-seq plasmids, termed "reporter RNAs",
109 are expressed at basal levels from the activity of the Ori itself. This allows detection of
110 silencing activity—the inhibition of the basal expression—in this assay.

111 Following its creation, the ATAC-STARR-seq plasmid library is transfected via
112 electroporation into a given cell line. From the same flask of cells, both reporter RNAs and
113 plasmid DNA are harvested 24 hours later, then prepared as Illumina sequencing libraries

114 and sequenced. Activity is calculated as the \log_2 ratio between normalized read counts from
115 the reporter RNA and plasmid DNA datasets. The re-isolation of plasmid DNA recovers only
116 the ATAC-STARR-seq plasmids that were successfully transfected, thus providing a more
117 accurate representation of the “input” sample than sequencing without transfection.
118 Supplemental Table S1 provides a comparison of experimental and analytical features as
119 well as reported data metrics for the current ATAC-STARR design and previously reported
120 approaches (Wang et al. 2018; Chaudhri et al. 2020).

121 **ATAC-STARR-seq maintains library complexity and nucleosome profiles of Tn5**
122 **selected DNA fragments**

123 Following the experimental design outlined above, we tagged GM12878 cells and
124 generated an ATAC-STARR-seq plasmid library that yielded about 50 million unique
125 accessible DNA fragments (Supplemental Text, Supplemental Fig. S1A). For a total of three
126 replicates, we then transfected the library into GM12878 cells and harvested both reporter
127 RNAs and plasmid DNA from the same flask of cells 24 hours later. We chose 24 hours
128 post-transfection to avoid significant effects from the plasmid-induced interferon gene
129 response, and to ensure the data reflects steady-state regulatory properties of GM12878
130 accessible regions (Supplemental Text, Supplemental Fig. S1B) (Muerdter et al. 2018).
131 Using the captured reporter RNAs and plasmid DNA, we prepared Illumina sequencing
132 libraries for each replicate and submitted for sequencing.

133 The size distribution of the accessible DNA fragments remained consistent throughout the
134 ATAC-STARR-seq procedure and displayed the characteristic nucleosome banding and
135 DNA pitch typified by ATAC-seq fragment libraries (Supplemental Fig. S2A-B). Analysis of
136 library complexity between replicates revealed an average maximum complexity of 90 million
137 unique fragments for input DNA, and 10 million unique fragments for reporter RNAs
138 (Supplemental Fig. S2C). The difference between RNA and DNA complexities is likely due to
139 higher duplication rates in the RNA samples (Supplemental Table S2) driven by both the

140 expression of multiple transcripts per plasmid and more PCR cycles required for the RNA
141 samples. In addition, for both RNA and DNA samples, replicates displayed high Pearson (r^2 :
142 0.96-0.99) and Spearman's (ρ : 0.77-0.93) correlation coefficients indicating strong
143 agreement among the three replicates assayed (Supplemental Fig. S3). Altogether the
144 ATAC-STARR-seq sequence libraries demonstrated the necessary quality and complexity
145 for downstream analysis.

146 **ATAC-STARR-seq faithfully captures chromatin accessibility with high signal-to-noise**

147 The use of Tn5 on native chromatin to selectively clone chromatin accessible DNA
148 fragments provides the opportunity to quantify not only reporter activity, but also chromatin
149 accessibility simultaneously from the same plasmid library. This is because the same DNA
150 fragments sequenced in a typical ATAC-seq workflow are contained in the ATAC-STARR-
151 seq plasmids. Given the insert fragments from reisolated plasmids are sequenced, we asked
152 if the resulting peak profiles recapitulate native ATAC-seq to measure chromatin accessibility.
153 This is important because, in contrast to a typical ATAC-seq procedure, ATAC-STARR-seq
154 involves several additional steps including cloning, transfection and reisolation, which could
155 distort the content of the library such that it no longer represents its native profile in the
156 genome. Specifically, mapped sequence reads derived from inserts of reisolated plasmids
157 are counted at a given locus and this estimate infers the accessibility of the region at the
158 time of tagmentation. This also reflects the number of plasmids that represent a given region
159 within the reisolated ATAC-STARR-seq plasmid library. To test this, we processed the
160 reisolated plasmid DNA as an Omni-ATAC-seq dataset and benchmarked against the
161 GM12878 Omni-ATAC-seq dataset from Corces *et al.* 2017. Raw sequences obtained for
162 both datasets were processed through identical workflows (see Methods). After collapsing
163 read duplicates, we called peaks for each dataset using a variety of false-discovery rates
164 (FDRs, Supplemental Table S3). To closely match the number of peaks previously reported
165 by Corces *et al.* 2017 (~108,433), we chose two separate FDR thresholds—0.0001 for

166 ATAC-STARR-seq and 0.001 for the Corces data—yielding 101,904 and 89,829 accessible
167 chromatin peaks respectively (Corces et al. 2017). The ATAC-STARR-seq and Corces *et al.*
168 peak sets represent 2.22% and 2.11% of the genome, respectively, which agrees with
169 previous reports (Figure 2A, (Thurman et al. 2012; Klemm et al. 2019)). Overall, 71% of
170 ATAC-STARR-seq peaks are reproduced in the Corces *et al.* dataset, while 81% of Corces
171 *et al.* peaks overlap the ATAC-STARR-seq dataset (Figure 2B; Jaccard index = 0.589),
172 indicating strong agreement between these data despite substantial differences in ATAC-
173 STARR DNA sample preparation. Furthermore, the fraction of reads in peaks score (FRiP),
174 an ENCODE ATAC-seq standard measure of noise, is considerably higher for both ATAC-
175 STARR-seq (0.74) and Corces *et al.* (0.526) than the ENCODE accepted standard (>0.2,
176 Figure 2C), indicating minimal background in our dataset. The high signal-to-noise is also
177 evident when looking at normalized read pileups at a representative locus (Figure 2D),
178 where the signal mirrors the Corces *et al.* accessibility signal patterns. Based on these
179 results, we conclude that ATAC-STARR-seq can accurately retain chromatin accessible
180 peaks in the human genome with high signal-to-noise.

181 **A sliding windows approach increases activity region calling sensitivity**

182 ATAC-STARR-seq tests regulatory activity in DNA enriched for accessible chromatin. Unlike
183 whole genome STARR-seq or other MPRAs, where the genomic DNA fragment distribution
184 is relatively constant, read coverage varies substantially from peak-to-peak in ATAC-
185 STARR-seq. In this way, ATAC-STARR-seq requires an analysis strategy that calls active
186 and silent regulatory regions within accessibility peaks. To address this “peaks-within-peaks”
187 problem, we developed an analytical approach using DESeq2 to normalize reporter RNA
188 read counts to reisolated plasmid DNA read counts. DESeq2 additionally performs an
189 independent filtering step which removes low count data confounders that can influence
190 ratios and result in false positive peak calls (Love et al. 2014).

191 We tested two different approaches for regulatory activity analysis. The two approaches
192 differ in how genomic regions are defined prior to differential analysis with DESeq2. Our
193 “sliding window” method, defines regions by slicing accessible peaks into 50bp sliding bins
194 with a 10bp step size (Figure 3A). Alternatively, the “fragment group” method, which is the
195 approach used in Wang *et al.* 2018, synthesizes regions by grouping paired-end sequencing
196 fragments by 75% or greater overlap (Supplemental Fig. S4A). Using a different set of
197 genomic regions, both methods assign and count overlapping RNA and DNA reads to each
198 genomic region and, using DESeq2, identify regions where the RNA count is statistically
199 different from the DNA count at a Benjamini-Hochberg (BH) adjusted p-value < 0.1 . The
200 “sliding window” method yielded ~30,000 distinct active regions, while the “fragment groups”
201 method yielded ~20,000 distinct active regions (Supplemental Fig. S4B). In addition, nearly
202 all active regions defined using the fragment group method (95%) are also captured in the
203 sliding window method regions (Supplemental Fig. S4C). Given this overlap and a 50%
204 greater recovery with the sliding windows approach, we used the sliding windows method to
205 call active ATAC-STARR-seq regulatory regions.

206 Because significance is the primary threshold in our region calling strategy, we examined the
207 influence of replicate count on the number of active regions called (Supplemental Text,
208 Supplemental Fig. S5). We found that, as expected, more replicates result in more active
209 regions. However, we caution that these additional regions may represent a disproportionate
210 number of false positives and may affect the outcomes of certain accuracy-sensitive
211 applications like computational modelling. In this way, we believe three replicates are
212 sufficient for most purposes. We also investigated the impact of read duplicates on replicate
213 correlations and activity region calling, finding that duplicate removal substantially hindered
214 region calling sensitivity despite yielding higher correlation coefficients between replicates
215 (described in more detail in Supplemental Text, Supplemental Fig. S6).

216 **Both short and long DNA fragments are required for comprehensive region calling**

217 Because DNA fragment synthesis for MPRAs is limited to 200 bp including the adapters and
218 barcode, a significant advantage of ATAC-STARR-seq and other capture-based MPRAs is
219 the ability to measure activity of longer DNA sequences (Santiago-Algarra et al. 2017). To
220 investigate the effect of fragment length on regulatory region calls, we divided mapped reads
221 into short (>125bp) and long (<125bp) fragments and independently called active and silent
222 regulatory regions; 125bp was chosen as it bisects the bimodal peak distribution displayed
223 by RNA and DNA libraries (Supplemental Fig. S2B). Overall, read counts were similar for
224 each sample after splitting into short and long groups (Supplemental Fig. S7A). Two to three
225 times as many active and silent regions were called in the long fragment group compared to
226 the short group (20,833 versus 10,789 for active and 16,872 versus 6,213 for silent).
227 Nonetheless, a substantial number of regions are called within the short fragment group,
228 although both fell short of the number of active and silent regions called when both long and
229 short were used (Supplemental Fig. S7B). The regulatory regions called using long DNA
230 fragments are larger than those called with short fragments, as expected (Supplemental Fig.
231 S7C); however, they display little difference in TSS distance, indicating these groups are not
232 comprised of different genomic annotations (Supplemental Fig. S7D). A critical observation
233 is that only 23% of active regions called using short reads overlap active regions called using
234 longer reads, revealing the two groups identify different regulatory regions in the genome
235 (Supplemental Fig. S7E); this is also true for the silent regulatory regions, although to a
236 lesser extent. Altogether this analysis reveals that short and long DNA fragments identify
237 different regulatory region sets both in number and similarity. Therefore, to be as
238 comprehensive as possible, STARR-seq assays should be designed to include both short
239 and long DNA fragments rather than impose a size selection to remove smaller fragments.

240 **ATAC-STARR-seq quantifies regulatory activity of open chromatin**

241 In the sliding window approach, bins are classified as active or silent depending on whether
242 RNA is enriched or depleted, respectively, and then like-bins are merged to collapse

243 overlaps (Figure 3A). Using this approach, we identified ~590,000 bins where RNA and DNA
244 counts were significantly different (Figure 3B). More specifically, this analysis identified
245 251,895 (4.1%) active bins and 339,737 (5.5%) silent bins from the ~5.6 million total bins
246 measured (Figure 3C). Overlapping bins were merged into 30,078 active and 21,125 silent
247 regulatory regions (Figure 3D). It is important to note that more silent than active bins are
248 called; however, because silent regions are generally larger (Figure 3E), merging
249 overlapping bins results in fewer silent regions than active. Collectively, the active and silent
250 bins represent ~9.5% of all bins measured, indicating that the majority of accessible DNA is
251 transcriptionally neutral. Moreover, most accessible peaks do not have an active or silent
252 region contained within them (69.5%), suggesting that most accessible regions are neutral
253 regulatory regions according to our assay (Figure 3F). This suggests that the majority of
254 accessible DNA has no regulatory potential in this cellular context or, alternatively, that
255 ATAC-STARR-seq is not sensitive enough to measure weakly active or weakly silent regions.
256 A recent study in mouse embryonic stem cells made the same observation using an
257 orthogonal approach, suggesting this phenomenon is present in other mammalian species
258 (Glaser et al. 2021). We note that a small percentage of accessible peaks (4.4%) contain
259 both active and silent regions, demonstrating that there can be competing regulatory regions
260 within the same accessible peak.

261 **Active and silent ATAC-STARR-seq regions represent both proximal and distal *cis*-**
262 **regulatory elements and lie within functional chromatin states**

263 To gain insight into the regulatory features of active regions, we annotated both active and
264 silent regions according to genomic location. Active regions are found in both promoter
265 proximal and distal areas of the genome, with a majority occurring in intronic and intergenic
266 sites (~55%), whereas silent regions coincide primarily with promoters (~75%) (Figure 4A).
267 Functional classification of active and silent regions by the 18-state ChromHMM model
268 (Roadmap Epigenomics Consortium et al. 2015) revealed that active regions consist of TSS

269 active, TSS flanking upstream, and Enhancer Active 1 chromatin states and are devoid of
270 repressive states like Repressed Polycomb Weak and Quiescent (Figure 4B). By contrast,
271 silent regions are slightly enriched for bivalent chromatin states (TSSBiv, EnhBiv), consistent
272 with the observation that they are accessible but not active. Most silent regions also coincide
273 with TSS Active and TSSFlank ChromHMM states, which corroborates their promoter
274 proximal locations; however, their designation as “active” by ChromHMM is somewhat
275 puzzling considering these DNA fragments do not drive transcription in our assay. One
276 explanation is that silent regulatory activity, as measured by episomal-based reporter assays,
277 does not fully copy regulatory activity as predicted by ChromHMM. Alternatively, active
278 promoters may confound the reporter assay by initiating transcription from the 3’UTR of the
279 plasmid causing conflicts with active transcription from the Ori.

280 To further investigate if silent regions are a result of 3’UTR transcription initiation, we
281 considered if an orientation bias existed in reporter RNAs levels. If 3’UTR transcription
282 conflicts exist, we would expect many fewer reporter RNAs when transcription results in
283 head-on conflicts rather than occurring in the same direction as the Ori. We therefore subset
284 reads based on whether they arose from an insert cloned in a 3’ to 5’ direction or in a 5’ to 3’
285 direction (Supplemental Fig. S8A). We then assigned read counts to all bins analyzed
286 (Supplemental Fig. S8B-C), the bins called active (Supplemental Fig. S8D-E), or the bins
287 called silent (Supplemental Fig. S8F-G). Because this is expected to be a promoter-specific
288 effect, we also split bins into proximal and distal based on location to the nearest
289 transcription start site. In all cases, more than 95% of the bins do not display an orientation
290 bias, which we defined as a normalized read count difference greater than five between
291 orientations (Supplemental Methods, Supplemental Fig. S8H). Moreover, we observe high
292 Pearson and Spearman’s correlation coefficients between orientations for all conditions (r^2 :
293 0.80-0.91 and p : 0.73-0.90) and the minimal contribution of orientation bias to silent regions
294 is in agreement with a previous report (Klein et al. 2020). For the <5% of regions that do
295 display orientation bias, proximal bins are more affected than distal bins, as expected.

296 Altogether, ATAC-STARR-seq does not display a significant orientation bias and most of the
297 21,000 silent regions we observe result from legitimate silencing activity or another source.

298 **Active and silent ATAC-STARR-seq regions are distinct functional classes and are**
299 **enriched for specific histone modifications and TF motifs**

300 To further investigate the chromatin landscape of the active and silent regions, we plotted
301 ENCODE GM12878 ChIP-seq signal (The ENCODE Project Consortium et al. 2020) for
302 EP300, CTCF, and histone modifications associated with active and repressed chromatin
303 states (Figure 4C). As expected, active regions contain EP300 at their center with histone 3
304 lysine 27 acetylation (H3K27ac) more broadly distributed across the center; histone 3 lysine
305 4 mono-methylation (H3K4me1) is also present at distal regions, while histone 3 lysine 4 tri-
306 methylation (H3K4me3) is at proximal regions. In addition, histone 3 lysine 27 tri-methylation
307 (H3K27me3)—a bivalent repressive mark—is largely absent from active regions. Proximal
308 silent regions, on the other hand, are enriched for H3K27me3 and H3K4me3. This suggests
309 many of the proximal silent regions are accessible bivalent regulatory elements in
310 lymphoblastoid cells. To support their designation as silent calls, we compared histone
311 modification signal at accessible peaks that contain either a silent region, an active region,
312 both a silent and active region, or neither, which we define as neutral accessible peaks
313 (Supplemental Fig. S9A). Consistent with the observations above, silent accessible peaks
314 contain more H3K27me3 signal and are devoid of H3K27ac signal relative to the other
315 accessible peak types.

316 It is important to note that silent regions are distinct from neutral regions, which are defined
317 as regions failing to reach significance in the RNA-DNA differential analysis. Overall, neutral
318 regions exhibit baseline levels of histone modifications and distribution in genomic
319 annotations like that of all accessible peaks (Supplemental Fig. S9B-C, Figure S4A-B). While
320 neutral regions represent the majority of accessible peaks, it is possible that a subset are

321 weak enhancers as indicated by overlap with ChromHMM states, or regulatory elements that
322 display activity in a different cellular context.

323 Our analysis of TF motifs within active and silent regions revealed prominent differences in
324 motif enrichment. Distal silent regions are strongly enriched for CTCF and its counterpart
325 BORIS, which is associated with diverse functions including gene repression and insulator
326 activity (Figure 4D-E)(Kim et al. 2015). In addition, we found enrichment for the SP/KLF
327 family several of which are known to be transcriptional repressors (Cao et al. 2010). By
328 contrast, the most enriched TFs in active regions were the IRF family, the ETS family,
329 subunits of the NF- κ B complex, and general promoter TFs such as THAP11 and YY1. These
330 data are consistent with our current understanding of immune gene regulation and regulatory
331 element function, which together corroborates the quantification of regulatory activity with
332 ATAC-STARR-seq.

333 **ATAC-STARR-seq retains the ability to map *in vivo* TF binding**

334 An inherent advantage of an ATAC-seq based approach is the ability to perform TF
335 footprinting. Computational footprinting methods identify Tn5 cleavage events or “cut sites”
336 from ATAC-seq data and, when combined with motif analysis, can identify TF binding sites
337 with high accuracy (Bentsen et al. 2020; Yan et al. 2020). Since ATAC-STARR-seq
338 produces similar high-quality chromatin accessibility peak profiles as standard ATAC-seq,
339 we explored whether TF footprints were also preserved. We generated Tn5-bias corrected
340 cut site signal files for both Corces *et al.* 2017 and ATAC-STARR-seq accessibility datasets
341 and plotted cut site signal at all accessible CTCF motifs (Figure 5A) (Bentsen et al. 2020).
342 As a control, we also plotted GM12878 CTCF ChIP-seq signal from ENCODE and ranked
343 region order by highest mean ChIP-seq signal. We observed consistent depletion of Tn5 cut-
344 sites for both Corces *et al.* 2017 and ATAC-STARR-seq accessibility at CTCF sites.
345 Moreover, we only observe footprints at motifs with CTCF ChIP-signal, demonstrating the
346 utility of TF footprinting to determine motifs that are bound or unbound by TFs. Given the

347 importance of TFs in driving enhancer function, this distinction is vital when dissecting
348 transcriptional regulation in human cells.

349 TF motif enrichment analysis pointed to multiple ETS family members, including ETS1 which
350 is an important immune cell regulator (Garrett-Sinha 2013) (Figure 4D). So, we asked
351 whether ETS1 footprints are also present in our data. Unlike CTCF, ETS1 shares its motif
352 with many other transcription factors, such as ETV4; therefore, footprinting cannot
353 distinguish ETS1 and ETV4 binding sites. For this reason, we refer to TFs using their
354 ENCODE-defined “archetypes”, which reflects the group of TFs that share the same motif
355 (Vierstra et al. 2020). For each archetype, we performed footprinting against one of the TFs
356 within an archetype to infer motifs bound by members of the group, such as ETS1 for the
357 ETS/1 archetype. To assess the extent to which ETS1 footprints can be explained by ETS1
358 binding, we plotted GM12878 ETS1 ChIP-seq signal from ENCODE within both Corces *et al.*
359 2017 and ATAC-STARR-seq cut sites (Figure 5B). Indeed, ETS1 ChIP-seq signal explains
360 the majority but not all the ETS/1 footprints present. We observe similar cut-site signal to
361 Corces *et al.* 2017, further indicating that ATAC-STARR-seq can detect *in vivo* binding of
362 transcription factors despite the additional cloning and transfection steps involved in
363 producing ATAC-STARR-seq DNA libraries.

364 We performed footprinting for several more immune related TF archetypes to identify bound
365 or unbound TF motifs (Figure 5C). For all TFs, bound motifs display substantially larger
366 footprint depth than unbound motifs. Together, this indicates that ATAC-STARR-seq, when
367 combined with footprinting, can identify regions of the genome where TFs are bound. This
368 additional level of information can be leveraged in conjunction with accessibility and activity
369 to understand the context of TF binding while circumventing the need to perform individual
370 chromatin immunoprecipitations.

371 **Collective profiling of accessibility, *in vivo* TF binding, and activity with ATAC-**
372 **STARR-seq reveals distinct networks of gene regulation**

373 Interrogating chromatin accessibility, TF binding, and regulatory activity together can be
374 used to interpret locus-specific gene regulatory mechanisms. For example, active regulatory
375 elements surrounding the B cell-specific expressed gene *ZBTB32* overlap IRF8 and NFKB1
376 footprints suggesting these regions are regulated by IRF8 and NFKB1 binding (Figure 6A).
377 We also observe SP1 and KLF3 footprints overlapping a silent region at the *ETV2* locus, a
378 gene lowly expressed in B cells, according to the Human Protein Atlas (Uhlen et al. 2015;
379 Uhlen et al. 2019). Together this indicates that active and silent regions can, in part, be
380 explained by the occupancy of these TFs.

381 To demonstrate the power of integrating TF footprints and regulatory regions on a global
382 scale, we clustered active and silent regions based on the presence or absence of several
383 TF footprints (Figure 6B-C). Footprints were selected based on top hits from the previous
384 motif enrichment analysis (Figure 4D-E). Regulatory activity may be driven by one or
385 multiple TF binding events that defines the cluster and is representative of a gene regulatory
386 network in the genome. Indeed, we find that the putative target genes regulated by each
387 unique group are enriched for distinct gene regulatory pathways and are often related to the
388 TFs in the cluster (Figure 6D-E). For example, cluster C is primarily defined by the presence
389 of IRF/1 and is enriched for interferon alpha/beta signalling. It is interesting that active
390 clusters tend to be more associated with B cell function than silent clusters, which are more
391 associated with general, non-B cell related pathways.

392 Altogether, these distinct gene regulatory networks provide an additional layer of insight into
393 the mechanisms that control gene expression and showcase how integration of the multiple
394 layers of gene regulatory information provided by ATAC-STARR-seq can narrow the focus of
395 gene targets for active and silent regions. We envision such an analysis could be used to
396 interpret the functional consequences of a dysregulated transcription factor or disease-
397 associated genetic variants. We provide this level of detail from a single dataset, which
398 further demonstrates the strong potential of our workflow to reveal distinct functional layers

399 of human gene regulation. The resolution we achieve here would not be possible without all
400 three levels of regulatory information provided by ATAC-STARR-seq.

401 **DISCUSSION**

402 Genome-wide approaches that integrate measurements of multiple layers of gene regulation
403 are needed to better understand enhancer function. By combining ATAC-seq with STARR-
404 seq, ATAC-STARR-seq assays regulatory activity only within the context of accessible
405 chromatin. This allows deeper coverage of regulatory elements by narrowing scope but
406 remaining inclusive of nearly all active regulatory elements. In this report, we substantially
407 expand the capabilities of ATAC-STARR-seq and present an improved workflow which
408 uniquely permits simultaneous profiling of accessibility, TF occupancy, and regulatory
409 activity from a single DNA fragment source. Specifically, we implement key experimental and
410 analytical improvements and present data rationalizing the decisions we make.

411 Experimentally, we adapt a modified tagmentation protocol (Omni-ATAC) to remove
412 mitochondrial DNA from the DNA fragment pool. We also utilize the Ori as the minimal
413 promoter on the STARR-seq backbone which improves reporter RNA expression, recovery,
414 and dynamic range over the super core promoter (SCP1) backbone (Muerdter et al. 2018;
415 Klein et al. 2020). Furthermore, we reisolate the transfected plasmid DNA to capture only the
416 DNA that is available to cells, which is a more accurate measure of the input than
417 sequencing prior to transfection. Reisolating plasmid DNA drives a greater degree of
418 variance between samples and better reflects a true experimental replicate than sequencing
419 the same DNA sample for each RNA replicate. Finally, we show that replicate number and
420 inclusion of long and short fragment sizes are critical for comprehensive region calling.

421 Critically, we developed and tested a simple and sensitive region calling strategy that
422 improves detection of regulatory regions including silencers. We also quantify chromatin
423 accessibility and identify TF footprints, which is surprising given the added processing steps
424 in ATAC-STARR-seq including cloning, transfection, and recapture of DNA libraries that can

425 dull or degrade footprint signal. This enabled us to stratify the active and silent regulatory
426 regions into distinct gene regulatory networks defined by the presence of one or multiple TF
427 footprints. Such an analysis typically requires multiple genomic sequencing assays, but we
428 do this using a single dataset.

429 With this improved workflow, we identified 30,078 active regions and 21,125 silent regions in
430 lymphoblastoid cells. Most active regions were distal to transcription start sites, enriched for
431 functional active ChromHMM states, and were enriched for known B cell regulating-TF
432 motifs such as IRF8 and NFkB. By contrast, the silencers are proximal to transcription start
433 sites and enriched for CTCF and the SP/KLF TF family. Silent regions are also enriched for
434 the bivalent marks H3K27me3 and H3K4me3 and may represent regulatory regions that are
435 poised, particularly at promoters. Because our plasmid design places regulatory regions
436 within the 3'UTR of the truncated reporter gene, it is possible that the lack of observed
437 reporter RNAs at silent regions are a result of head-on transcriptional conflicts that arise
438 from antisense transcription initiation from the 3'UTR. However, we show this minimally
439 occurs in our system and the silent regions reflect true silencing activity or another source
440 that has yet to be identified. While further studies may be needed to validate these silent
441 regions, this work confirms that the silencers are a distinct class of regulatory element with
442 distinct properties compared to active and neutral regions and warrant further investigation.
443 Even with an increasing number of studies targeted at identifying silencers in the human
444 genome, silencing regulatory regions remain an under-studied aspect of gene regulation and
445 our approach provides a new strategy to investigate these elements on a global scale (Doni
446 Jayavelu et al. 2020; Pang and Snyder 2020; Kim et al. 2021).

447 ATAC-STARR-seq data has several distinct attributes that require a tailored analysis
448 strategy. Current MPRA bioinformatic tools and pipelines are not tractable for these data,
449 because in ATAC-STARR-seq the input itself is enriched for accessible chromatin and the
450 read pileup varies considerably within these loci. In this way, the analysis of our data

451 required calling essentially “peaks within peaks”. For this reason, it was critical to 1)
452 normalize RNA to DNA and 2) avoid regions of low count data, which is why we adapted
453 approaches using DESeq2. We also showed that including PCR duplicates was preferred
454 over collapsing duplicates. In the future it would be beneficial to introduce a unique
455 molecular identifier to the system—such as the strategy employed by UMI-STARR-seq
456 (Neumayr et al. 2019)—to collapse only the duplicates arising from PCR. While we show
457 comparisons of analysis strategies here, we believe that more information could be extracted
458 from this and future ATAC-STARR-seq datasets with improved analysis strategies. In recent
459 years we have seen the development of tailormade peak callers for whole genome STARR-
460 seq, such as CRADLE (Kim et al. 2021) and STARRPeaker (Lee et al. 2020); a similarly
461 tailored ATAC-STARR-seq peak caller could further improve the capabilities of the method.

462 While this study was limited to one condition, there are many potential applications of ATAC-
463 STARR-seq. With the ability to subset enhancers by TF occupancy, ATAC-STARR-seq
464 could be leveraged to investigate enhancer grammar by pairing measurable regulatory
465 activity with multiple TF footprints that drive enhancer function. This approach also has the
466 potential to identify dysfunctional gene regulatory networks in diseases like cancer where
467 neoplastic transformation can be driven by the dysfunction of a specific TF. Additionally, an
468 ATAC-STARR-seq plasmid library may be generated from one cell-type and tested in
469 another. This flexibility could be used as a tool to determine context dependent activity or
470 investigate enhancer and TF usage patterns during a differentiation time course.

471 In this study, we demonstrated that our improved ATAC-STARR-seq workflow is a powerful
472 approach enabling joint quantification of chromatin accessibility, transcription factor
473 occupancy, and regulatory activity. We further demonstrate how this single assay can
474 characterize the human genome at many functional levels from chromatin accessibility to
475 distinct gene regulatory networks. This method provides a state-of-the-art approach to
476 deeply investigate transcriptional regulation of the human genome. We provide a detailed

477 protocol and a well-documented code repository so that ATAC-STARR-seq may be easily
478 used and adapted by the field.

479 **MATERIAL AND METHODS**

480 **Cell Culture**

481 GM12878 cells were obtained from Coriell and cultured with RPMI 1640 Media containing 15%
482 fetal bovine serum, 2mM GlutaMAX, 100 units/mL penicillin and 100 µg/mL streptomycin.
483 Cells were cultured at 37°C, 80% relative humidity, and 5% CO₂. Cell density was
484 maintained between 0.2×10⁶ and 1×10⁶ cells/mL with a 50% media change every 2-4 days.
485 All cell lines were regularly screened for mycoplasma contamination using the MycoAlert kit
486 (Lonza).

487 **Plasmids**

488 The hSTARR-seq_ORI plasmid vector was a gift from Alexander Stark (Addgene plasmid
489 #99296) and the pcDNA3-EGFP was a gift from Doug Golenbock (Addgene plasmid
490 #13031). The bacterial stabs from Addgene were spread onto an LB agar plate containing
491 100µg/mL ampicillin and incubated at 37°C overnight. For each, a single colony was picked
492 and grown in 50mL LB containing 100µg/mL ampicillin overnight at 37°C while shaking at
493 225rpm. Plasmid DNA was extracted using the ZymoPURE II Plasmid Midiprep kit (Zymo
494 Research, #D4200).

495 The linear vector used in the ATAC-STARR-seq gibson cloning step was generated by a
496 single 50µL PCR reaction using NEBNext® Ultra™ II Q5® Master Mix (NEB, #M0544S).
497 While not necessary for this study, primers were designed to add the i5 barcode to the
498 linearized vector; this allows for different ATAC-STARR-seq plasmid libraries to be pooled
499 and tracked. Following this approach, a universal forward primer (Fwd_universal_STARR)
500 and a reverse primer (Rev_N504_STARR) designed to add the N504 barcode were used

501 (primer sequences are provided in Supplemental Table S4). PCR Products were purified
502 with the Zymo Research DNA Clean & Concentrator-5 kit. DNA yield was determined by
503 Nanodrop, and purity was analysed by gel electrophoresis; the linearized vector was the only
504 product observed on the gel.

505 **Tagmentation**

506 A total of eight tagmentation reactions were performed on 50,000 GM12878 cells each. We
507 followed a slightly modified version of the Omni-ATAC approach used in Corces *et al.* 2017
508 (Corces *et al.* 2017). Specifically, twice as much Tn5 than described in the protocol was
509 used. Standard Tn5 transposase was prepared in-house following the method described in
510 Picelli *et al.* 2014 (Picelli *et al.* 2014). Standard Tn5 transposome was assembled as
511 described in Barnett *et al.* 2020 (Barnett *et al.* 2020) with the following oligos: Tn5_1,
512 Tn5_2_ME_comp, and TN5MEREV. Tagmented products were pooled together and purified
513 with the Zymo Research DNA Clean & Concentrator-5 kit (#D4013). The entire elution was
514 split and amplified via five-10 μ L PCR reactions. We used NEBNext® High-Fidelity 2 \times PCR
515 Master Mix (#M0541S)—which is not a hot-start formulation—to first extend tagments before
516 the initial denaturation step of PCR via the following cycling parameters: 72°C 5 min, 98°C
517 30s; 4 cycles of 98°C 10s, 62°C 30s, 72°C 60s; final extension 72°C 2 min; hold at 10°C.
518 Forward and reverse primer sequences, Fwd_atac-starr_tag and Rev_atac-starr_tag, are
519 provided in Supplemental Table S3. Amplified products were purified with the Zymo
520 Research DNA Clean & Concentrator-5 kit and then analyzed for concentration and size
521 distribution with a HSD5000 screentape (Agilent, #5067) on an Agilent 4150 TapeStation
522 system. After amplification, we selected PCR products less than 500bp using SPRISelect
523 beads (Beckman-Coulter, #B23317) at a 0.6 \times volume ratio of beads:sample. Selection was
524 verified using a HSD5000 screentape.

525 **Massively Parallel Cloning**

526 Four 10 μ L gibson cloning reactions were performed with NEBuilder® HiFi DNA Assembly
527 Master Mix at a vector:insert molar ratio of 1:2. As a negative control, we performed one
528 cloning reaction substituting tagments with nuclease-free water. Gibson products were
529 pooled and purified via ethanol precipitation as previously described in Sambrook & Russell
530 (Sambrook and Russell 2006); we used glycoblu (150 μ g/mL) as a co-precipitant. Purified
531 gibson products were electroporated into MegaX DH10B T1R Electrocomp™ Cells
532 (Invitrogen, # C640003) using a Bio-Rad Gene Pulser. Three electroporations for the ATAC-
533 STARR-seq sample (and 1 for the control) were performed with the following parameters:
534 exponential decay pulse type, 2kV, 200 Ω , 25 μ F, and 0.1cm gap distance. Pre-warmed SOC
535 media (1mL) was added immediately following electroporation; the three reactions were
536 pooled and incubated at 37°C for 1 hour. We confirmed cloning success by plating a dilution
537 series—using a small aliquot from the ATAC-STARR-seq and negative control samples—
538 onto pre-warmed LB agar plates containing 100 μ g/mL ampicillin and visualizing colonies 24
539 hours later. The remaining ATAC-STARR-seq transformation was added directly to a 1L LB
540 liquid culture with 100 μ g/mL ampicillin and grown at 37°C while shaking at 225rpm overnight.
541 The next day, plasmid DNA was harvested from the 1L culture using the ZymoPURE II
542 Plasmid Gigaprep (Zymo Research, #D4204). Before prepping, we recorded a 1.633 optical
543 density.

544 **Electroporation**

545 GM12878 cells were cultured so that cell density was between 400,000 and 800,000
546 cells/mL on day of transfection. Three replicates were performed on separate days. For each
547 replicate, a total of 20 electroporation reactions was performed using the Neon™
548 Transfection System 100 μ L Kit (Invitrogen, #MPK10025) and the associated Neon™
549 Transfection System (Invitrogen, #MPK5000). 121 million GM12878 cells were collected,
550 washed with 45mL PBS, and resuspended in 2178 μ L Buffer R. For each reaction, 5 million
551 cells (in 90 μ L Buffer R) were electroporated with 5 μ g of ATAC-STARR-seq plasmid DNA (in

552 10 μ L nuclease-free water) in a total volume of 100 μ L with the following parameters: 1100V,
553 30ms, and 2 pulses. Electroporated cells were dispensed immediately into a pre-warmed T-
554 75 flask containing 50mL of RPMI 1640 with 20% fetal bovine serum and 2mM GlutaMAX.

555 **Cell Harvest**

556 24 hours after transfection, the 50mL ATAC-STARR-seq flask was divided into two equal
557 volumes; plasmid DNA was extracted from one volume, while reporter RNAs were extracted
558 from the other. Plasmid DNA was isolated with the ZymoPURE II Plasmid Midiprep kit
559 (#D4200) and eluted in 50 μ L 10mM Tris-HCL pH 8.0. Prior to lysis, cells were washed with
560 25mL PBS to remove any extracellular plasmid DNA. Reporter RNAs were extracted in a
561 stepwise process. First, total RNA was isolated from the second volume of transfected cells
562 using the TRIzol™ Reagent and Phasemaker™ Tubes Complete System (Invitrogen™,
563 #A33251). Specifically, 5mL TRIzol was added to homogenize the washed and pelleted cells.
564 Next, polyadenylated RNA was isolated from total RNA using oligo(dT)25 Magnetic Beads
565 (NEB, #S1419S) at a 1 μ g Total RNA to 10 μ g beads ratio. We performed this step at 4°C and
566 eluted into 50 μ L 10mM Tris-HCl pH 7.5. The extracted poly(A)⁺ RNA was treated with DNase
567 I (NEB, #M0303S). This reaction was cleaned up using the Zymo Research RNA Clean &
568 Concentrator-25 kit (Zymo Research, #R1018).

569 **First-strand cDNA synthesis**

570 For each sample, ten 50 μ L reverse transcription reactions were carried out using
571 PrimeScript™ Reverse Transcriptase (Takara, #2680) and a gene specific primer
572 (STARR_GSP) as described by Muerdter *et al.* 2018 (Muerdter et al. 2018). Single-stranded
573 cDNA was treated with RNase A at a concentration of 20 μ g/mL in low salt concentrations
574 and cleaned up with a Zymo Research DNA Clean & Concentrator-5 kit.

575 **Illumina Sequencing Library Preparation**

576 For reisolated plasmid and reporter RNA samples, Illumina-compatible libraries were
577 generated using NEBNext® Ultra™ II Q5® Master Mix and a unique combination of the
578 following Nextera indexes: N504-N505 (i5) and N701-N702 (i7), see Supplemental Table S1
579 for primer sequences. DNA samples were amplified for 8 PCR cycles, while RNA was
580 amplified for 12-13 cycles. In both cases, products were purified with the Zymo Research
581 DNA Clean & Concentrator-5 kit and analyzed for concentration and size distribution using a
582 HSD5000 screentape. Purified products were sequenced on an Illumina NovaSeq, PE150,
583 at a requested read depth of 50 or 75 million reads, for DNA and RNA samples, respectively,
584 on an Illumina NovaSeq 6000 machine through the Vanderbilt Technology for Advanced
585 Genomics (VANTAGE) sequencing core. Reads were processed and analyzed as described
586 in the supplemental methods. We provide guidelines for ATAC-STAR-seq quality control in
587 the supplemental text.

588 **DATA ACCESS**

589 All raw and processed sequencing data generated in this study have been submitted to the
590 NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession
591 number GSE181317. Python scripts and additional code for ATAC-STARR-seq data
592 analysis are available at GitHub ([https://github.com/HodgesGenomicsLab/ATAC-STARR-](https://github.com/HodgesGenomicsLab/ATAC-STARR-seq)
593 [seq](https://github.com/HodgesGenomicsLab/ATAC-STARR-seq)) and as Supplemental Code. An interactive version of the protocol is posted
594 on [protocols.io](https://www.protocols.io) ([dx.doi.org/10.17504/protocols.io.b2nugdew](https://doi.org/10.17504/protocols.io.b2nugdew)) and a pdf version of the
595 protocol at publication date is included as a Supplemental file.

596 **COMPETING INTEREST STATEMENT**

597 None declared.

598 **ACKNOWLEDGEMENTS**

599 We thank Felix Muerdter, Sarah Fong, Tony Capra, and members of the Hodges Lab,
600 especially Kelly Barnett, Tim Scott, Lindsey Guerin, Verda Agan, Elizabeth Dorans, and Ali
601 Wilt for helpful feedback and discussions. We also thank Biorender.com for illustrations,
602 Addgene for plasmids used in the study, the Dave Cortez Lab for use of their Bio-Rad Gene
603 Pulser, and the Manny Ascano Lab for qPCR primers and helpful advice. We are grateful for
604 support of the project and the time invested in producing this manuscript by the NIH awards
605 [K22 CA184308-03 to E.H], Department of Defense Idea Award [W81XWH-20-1-0522 to
606 E.H], and American Cancer Society (ACS) Institutional Research Grant (#IRG-15-169-56).

607 REFERENCES

- 608 Arnold CD, Gerlach D, Stelzer C, Boryn LM, Rath M, Stark A. 2013. Genome-wide
609 quantitative enhancer activity maps identified by STARR-seq. *Science* **339**: 1074-
610 1077.
- 611 Barnett KR, Decato BE, Scott TJ, Hansen TJ, Chen B, Attalla J, Smith AD, Hodges E. 2020.
612 ATAC-Me Captures Prolonged DNA Methylation of Dynamic Chromatin Accessibility
613 Loci during Cell Fate Transitions. *Mol Cell* **77**: 1350-1364 e1356.
- 614 Bentsen M, Goymann P, Schultheis H, Klee K, Petrova A, Wiegandt R, Fust A, Preussner J,
615 Kuenne C, Braun T et al. 2020. ATAC-seq footprinting unravels kinetics of
616 transcription factor binding during zygotic genome activation. *Nat Commun* **11**: 4267.
- 617 Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. 2013. Transposition of native
618 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding
619 proteins and nucleosome position. *Nat Methods* **10**: 1213-1218.
- 620 Cao Z, Sun X, Icli B, Wara AK, Feinberg MW. 2010. Role of Kruppel-like factors in leukocyte
621 development, function, and disease. *Blood* **116**: 4404-4414.
- 622 Chaudhri VK, Dienger-Stambaugh K, Wu Z, Shrestha M, Singh H. 2020. Charting the cis-
623 regulome of activated B cells by coupling structural and functional genomics. *Nat*
624 *Immunol* **21**: 210-220.
- 625 Chen AF, Parks B, Kathiria AS, Ober-Reynolds B, Goronzy JJ, Greenleaf WJ. 2022. NEAT-
626 seq: simultaneous profiling of intra-nuclear proteins, chromatin accessibility and gene
627 expression in single cells. *Nat Methods* **19**: 547-553.
- 628 Clark SJ, Argelaguet R, Kapourani CA, Stubbs TM, Lee HJ, Alda-Catalinas C, Krueger F,
629 Sanguinetti G, Kelsey G, Marioni JC et al. 2018. scNMT-seq enables joint profiling of
630 chromatin accessibility DNA methylation and transcription in single cells. *Nat*
631 *Commun* **9**: 781.
- 632 Corces MR, Trevino AE, Hamilton EG, Greenside PG, Sinnott-Armstrong NA, Vesuna S,
633 Satpathy AT, Rubin AJ, Montine KS, Wu B et al. 2017. An improved ATAC-seq
634 protocol reduces background and enables interrogation of frozen tissues. *Nat*
635 *Methods* **14**: 959-962.
- 636 Doni Jayavelu N, Jajodia A, Mishra A, Hawkins RD. 2020. Candidate silencer elements for
637 the human and mouse genomes. *Nat Commun* **11**: 1061.
- 638 Garrett-Sinha LA. 2013. Review of Ets1 structure, function, and roles in immunity. *Cell Mol*
639 *Life Sci* **70**: 3375-3390.
- 640 Gasperini M, Tome JM, Shendure J. 2020. Towards a comprehensive catalogue of validated
641 and target-linked human enhancers. *Nat Rev Genet* **21**: 292-310.

- 642 Glaser LV, Steiger M, Fuchs A, van Bommel A, Einfeldt E, Chung HR, Vingron M, Meijsing
643 SH. 2021. Assessing genome-wide dynamic changes in enhancer activity during
644 early mESC differentiation by FAIRE-STARR-seq. *Nucleic Acids Res* **49**: 12178-
645 12195.
- 646 Heinz S, Romanoski CE, Benner C, Glass CK. 2015. The selection and function of cell type-
647 specific enhancers. *Nat Rev Mol Cell Biol* **16**: 144-154.
- 648 Inoue F, Kircher M, Martin B, Cooper GM, Witten DM, McManus MT, Ahituv N, Shendure J.
649 2017. A systematic comparison reveals substantial differences in chromosomal
650 versus episomal encoding of enhancer activity. *Genome Res* **27**: 38-52.
- 651 Johnson GD, Barrera A, McDowell IC, D'Ippolito AM, Majoros WH, Vockley CM, Wang X,
652 Allen AS, Reddy TE. 2018. Human genome-wide measurement of drug-responsive
653 regulatory activity. *Nat Commun* **9**: 5317.
- 654 Kelly TK, Liu Y, Lay FD, Liang G, Berman BP, Jones PA. 2012. Genome-wide mapping of
655 nucleosome positioning and DNA methylation within individual DNA molecules.
656 *Genome Res* **22**: 2497-2506.
- 657 Kim S, Yu NK, Kaang BK. 2015. CTCF as a multifunctional protein in genome regulation and
658 gene expression. *Exp Mol Med* **47**: e166.
- 659 Kim YS, Johnson GD, Seo J, Barrera A, Cowart TN, Majoros WH, Ochoa A, Allen AS,
660 Reddy TE. 2021. Correcting signal biases and detecting regulatory elements in
661 STARR-seq data. *Genome Res* doi:10.1101/gr.269209.120.
- 662 Kircher M, Xiong C, Martin B, Schubach M, Inoue F, Bell RJA, Costello JF, Shendure J,
663 Ahituv N. 2019. Saturation mutagenesis of twenty disease-associated regulatory
664 elements at single base-pair resolution. *Nat Commun* **10**: 3583.
- 665 Klein JC, Agarwal V, Inoue F, Keith A, Martin B, Kircher M, Ahituv N, Shendure J. 2020. A
666 systematic evaluation of the design and context dependencies of massively parallel
667 reporter assays. *Nat Methods* **17**: 1083-1091.
- 668 Klemm SL, Shipony Z, Greenleaf WJ. 2019. Chromatin accessibility and the regulatory
669 epigenome. *Nat Rev Genet* **20**: 207-220.
- 670 Lee D, Shi M, Moran J, Wall M, Zhang J, Liu J, Fitzgerald D, Kyono Y, Ma L, White KP et al.
671 2020. STARRPeaker: uniform processing and accurate identification of STARR-seq
672 active regions. *Genome Biol* **21**: 298.
- 673 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
674 RNA-seq data with DESeq2. *Genome Biol* **15**: 550.
- 675 Maricque BB, Dougherty JD, Cohen BA. 2017. A genome-integrated massively parallel
676 reporter assay reveals DNA sequence determinants of cis-regulatory activity in
677 neural cells. *Nucleic Acids Res* **45**: e16.
- 678 Melnikov A, Murugan A, Zhang X, Tesileanu T, Wang L, Rogov P, Feizi S, Gnirke A, Callan
679 CG, Jr., Kinney JB et al. 2012. Systematic dissection and optimization of inducible
680 enhancers in human cells using a massively parallel reporter assay. *Nat Biotechnol*
681 **30**: 271-277.
- 682 Muerdt F, Boryn LM, Woodfin AR, Neumayr C, Rath M, Zabidi MA, Pagani M, Haberle V,
683 Kazmar T, Catarino RR et al. 2018. Resolving systematic errors in widely used
684 enhancer activity assays in human cells. *Nat Methods* **15**: 141-149.
- 685 Neumayr C, Pagani M, Stark A, Arnold CD. 2019. STARR-seq and UMI-STARR-seq:
686 Assessing Enhancer Activities for Genome-Wide-, High-, and Low-Complexity
687 Candidate Libraries. *Curr Protoc Mol Biol* **128**: e105.
- 688 Pang B, Snyder MP. 2020. Systematic identification of silencers in human cells. *Nat Genet*
689 **52**: 254-263.
- 690 Patwardhan RP, Hiatt JB, Witten DM, Kim MJ, Smith RP, May D, Lee C, Andrie JM, Lee SI,
691 Cooper GM et al. 2012. Massively parallel functional dissection of mammalian
692 enhancers in vivo. *Nat Biotechnol* **30**: 265-270.
- 693 Picelli S, Bjorklund AK, Reinius B, Sagasser S, Winberg G, Sandberg R. 2014. Tn5
694 transposase and tagmentation procedures for massively scaled sequencing projects.
695 *Genome Res* **24**: 2033-2040.

- 696 Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,
697 Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J et al. 2015. Integrative analysis
698 of 111 reference human epigenomes. *Nature* **518**: 317-330.
- 699 Sambrook J, Russell DW. 2006. Standard ethanol precipitation of DNA in microcentrifuge
700 tubes. *CSH Protoc* **2006**.
- 701 Santiago-Algarra D, Dao LTM, Pradel L, Espana A, Spicuglia S. 2017. Recent advances in
702 high-throughput approaches to dissect enhancer function. *F1000Res* **6**: 939.
- 703 The ENCODE Project Consortium, Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shores N,
704 Adrian J, Kawli T, Davis CA, Dobin A et al. 2020. Expanded encyclopaedias of DNA
705 elements in the human and mouse genomes. *Nature* **583**: 699-710.
- 706 Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, Sheffield NC,
707 Stergachis AB, Wang H, Vernot B et al. 2012. The accessible chromatin landscape of
708 the human genome. *Nature* **489**: 75-82.
- 709 Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, Sivertsson A,
710 Kampf C, Sjostedt E, Asplund A et al. 2015. Tissue-based map of the human
711 proteome. *Science* **347**: 1260419.
- 712 Uhlen M, Karlsson MJ, Zhong W, Tebani A, Pou C, Mikes J, Lakshmikanth T, Forsstrom B,
713 Edfors F, Odeberg J et al. 2019. A genome-wide transcriptomic analysis of protein-
714 coding genes in human blood cells. *Science* **366**.
- 715 Vierstra J, Lazar J, Sandstrom R, Halow J, Lee K, Bates D, Diegel M, Dunn D, Neri F,
716 Haugen E et al. 2020. Global reference mapping of human transcription factor
717 footprints. *Nature* **583**: 729-736.
- 718 Wang X, He L, Goggin SM, Saadat A, Wang L, Sinnott-Armstrong N, Claussnitzer M, Kellis
719 M. 2018. High-resolution genome-wide functional dissection of transcriptional
720 regulatory regions and nucleotides in human. *Nat Commun* **9**: 5380.
- 721 Yan F, Powell DR, Curtis DJ, Wong NC. 2020. From reads to insight: a hitchhiker's guide to
722 ATAC-seq data analysis. *Genome Biol* **21**: 22.

723

724 **FIGURE LEGENDS**

725 **Figure 1. Schematic of the ATAC-STARR-seq methodology.** (A) The experimental design
726 of ATAC-STARR-seq consists of three parts: plasmid library generation, reporter assay, and
727 data analysis. Open chromatin is isolated from cells with the cut and paste transposase Tn5
728 and only large DNA fragments (>500bp) are removed. The open chromatin fragments are
729 cloned into a reporter plasmid and the resulting clones—called an ATAC-STARR-seq
730 plasmid library—are electroporated into cells. 24 hours later, both reporter RNAs (blue)—
731 which are transcribed directly off the ATAC-STARR-seq plasmid—and ATAC-STARR-seq
732 plasmid DNA (red) are harvested, and Illumina-sequencing libraries are prepared and
733 sequenced. The resulting ATAC-STARR-seq sequence data is analyzed to extract
734 regulatory activity, chromatin accessibility, and transcription factor footprints. (B) Reporter
735 plasmid design and the expected outcomes for neutral, active, and silent regulatory

736 elements. Each ATAC-STARR-seq plasmid within a library contains a truncated GFP (trGFP)
737 coding sequence, a poly-adenylation signal sequence, an origin of replication (Ori) (which
738 moonlights as a minimal core promoter), and the unique open chromatin fragment being
739 assayed. Since the accessible region is contained in the 3' UTR, the abundance of itself in
740 the transcript pool reflects its activity. In this way, neutral elements do not affect the system
741 and reporter RNAs are expressed at a basal expression level dictated by the minimal core
742 promoter, the Ori. Accessible chromatin fragments that are active express reporter RNAs at
743 a higher level than the basal expression level, while silent elements repress the Ori and
744 reporter RNAs are expressed at a lower level than basal expression. Dashed boxes
745 represent new components of the ATAC-STARR-seq assay design and workflow.

746 **Figure 2. ATAC-STARR-seq accurately quantifies chromatin accessibility.** ATAC-seq
747 data from Corces et al. 2017 is compared with ATAC-STARR-seq plasmid DNA data. (A)
748 Fraction of the human genome represented by each peak set. (B) Venn diagram of peak
749 overlap between the two datasets and the associated Jaccard Index. (C) Fraction of paired-
750 end (PE) fragments in peaks—FRiP scores—for both samples. (D) Signal tracks comparing
751 counts per million (CPM) normalized read count at a representative locus.

752 **Figure 3. ATAC-STARR-seq quantifies regulatory activity within accessible chromatin.**
753 (A) Schematic of the sliding window peak calling method. Accessibility peaks are chopped
754 into 50bp bins at a 10bp step size with the BEDTools makewindows function (options -w 50,
755 -s 10). For each window, RNA and DNA reads are counted using Subread's featureCounts
756 function. Differential analysis comparing RNA and DNA read count is performed with
757 DESeq2. Significant bins are called at an Benjamini-Hochberg (BH) adjusted p-value < 0.1
758 and parsed into active or silent depending on \log_2 fold-change (FC) value (+/- zero). Finally,
759 bins are collapsed into regions using the BEDTools merge function. \log_2 FC scores are
760 averaged across merged bins. (B) Volcano plot of \log_2 FC scores against $-\log_{10}$ -transformed
761 BH adjusted p-value from DESeq2 for all bins analyzed. (C) The proportion of bins called as
762 active or silent. (D) The number of regions defined as either active or silent. (E) Overlapping

763 density plots of active and silent regulatory region size; dashed lines represent the medians
764 in each case. (F) The proportion of accessible peaks that overlap an active or silent region,
765 or both.

766 **Figure 4. Regulatory regions defined by ATAC-STARR exhibit annotations, histone**
767 **modifications, and TFs characteristic of their function.** (A) Annotation of regulatory
768 regions relative to the transcriptional start site (TSS). The promoter is defined as 2kb
769 upstream and 1kb down-stream of the TSS. (B) Annotation of regulatory regions by the
770 ChromHMM 18-state model for GM12878 cells. (C) Heatmaps of GM12878 ENCODE ChIP-
771 seq signal and regulatory activity for proximal and distal ATAC-STARR-defined regulatory
772 regions. Proximal regions were classified as within 2kb upstream and 1kb downstream of a
773 TSS; all other regions were annotated as distal. Active and silent regions were ranked by
774 mean activity signal for both proximal and distal regions. (D-E) Transcription factor motif
775 enrichment analysis as quantified by HOMER. Fold-change values are relative to the default
776 background calculated by HOMER.

777 **Figure 5. ATAC-STARR-seq identifies transcription factor footprints.** (A) Comparison of
778 ENCODE CTCF ChIP-seq signal to Corces et al. and ATAC-STARR-seq cut count signal for
779 all accessible CTCF motifs. (B) Comparison of ENCODE ETS1 ChIP-seq signal to Corces et
780 al. and ATAC-STARR-seq cut count signal for all accessible motifs with the ETS/1 motif
781 archetype. For both, regions were ranked by largest mean ChIP-seq signal. (C) Aggregate
782 plots representing mean signal for the TOBIAS-defined bound and unbound motif
783 archetypes: CTCF, ETS/1, CREB/ATF/1, IRF/1, SPI, NFKB/2.

784 **Figure 6. TF footprints stratify ATAC-STARR-defined regulatory regions into gene**
785 **regulatory networks.** (A) ATAC-STARR-defined chromatin accessibility, TF footprints, and
786 regulatory regions at Chr19:35,611,232-35,798,446 (hg38). Signal tracks represent counts
787 per million normalized read depth of chromatin accessibility. Zooms into *ETV2* and *ZBTB32*
788 show that some regulatory regions are occupied by a SP1, KLF3, IRF8, or NFKB1 footprint.

789 (B-C) Heatmaps of clustered (B) active and (C) silent regions based on presence or absence
790 of footprints for select TF motif archetypes. (D-E) Reactome pathway enrichment analysis for
791 nearest-neighbor gene sets for each of the clusters. Genes counts for each cluster are
792 displayed below their group identifier.









