

1 **Computational Analysis of Maize Enhancer Regulatory Elements Using ATAC-STARR-seq**

2 Alexandre P. Marand*

3
4 Department of Genetics, University of Georgia, Athens, USA

5 *For correspondence: marand@uga.edu

6
7 **[Abstract]** The blueprints to development, response to the environment, and cellular function are
8 largely the manifestation of distinct gene expression programs controlled by the spatiotemporal activity
9 of *cis*-regulatory elements. Although biochemical methods for identifying accessible chromatin – a
10 hallmark of active *cis*-regulatory elements – have been developed, approaches capable of measuring
11 and quantifying *cis*-regulatory activity are only beginning to be realized. Massively Parallel Reporter
12 Assays coupled to chromatin accessibility profiling present a high-throughput solution for testing the
13 transcription-activating capacity of millions of putatively regulatory DNA sequences in parallel.
14 However, clear computational pipelines for analyzing these high-throughput sequencing-based
15 reporter assays are lacking. In this protocol, I layout and rationalize a computational framework for the
16 processing and analysis of Assay for Transposase Accessible Chromatin profiling followed by Self-
17 Transcribed Active Regulatory Region sequencing (ATAC-STARR-seq) data from a recent study in *Zea*
18 *mays*. The approach described herein can be adapted to other sequencing-based reporter assays and
19 is largely agnostic to the model organism with the appropriate input substitutions.

20
21 **Keywords:** STARR-seq, ATAC-seq, ATAC-STARR-seq, Regulatory activity, *cis*-regulatory elements,
22 Accessible chromatin, Transcriptional regulation, Enhancers

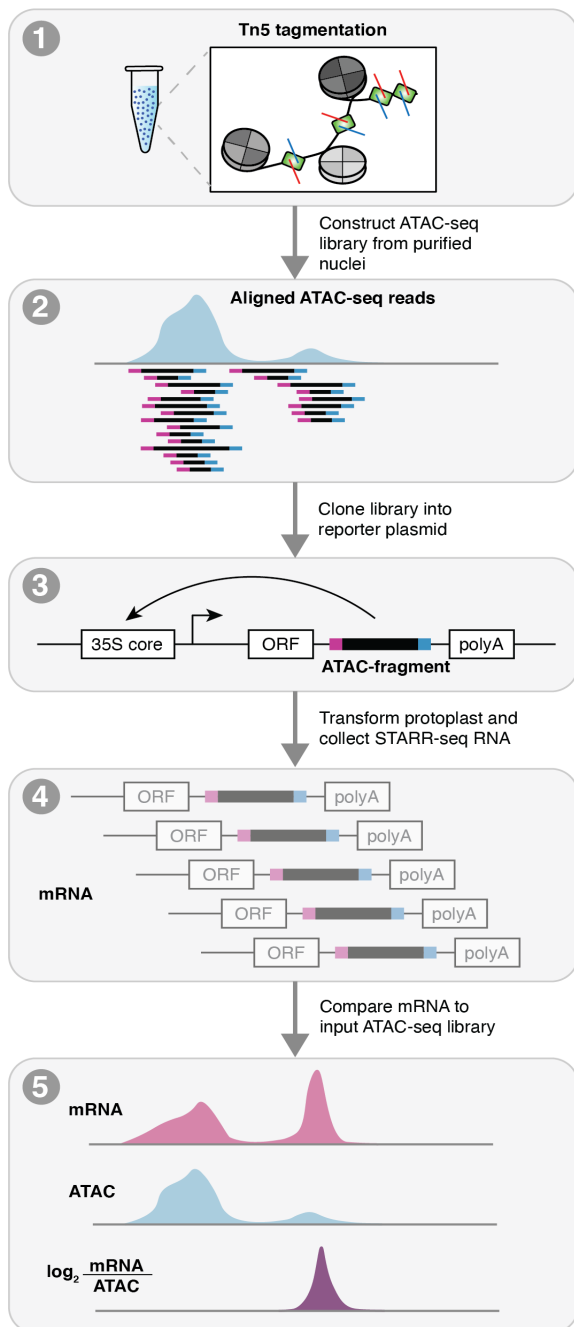
23
24 **[Background]** Eukaryotic cells exhibit remarkable functional and morphological diversity despite
25 containing a generally invariant copy of the same genomic sequence. Cellular heterogeneity arises in
26 part due to the activities of *cis*-regulatory elements (CREs), short DNA binding motifs recognized by
27 sequence specific transcription factors (TFs). CREs are often found in clusters termed *cis*-regulatory
28 modules (CRMs) that dictate highly dynamic spatiotemporal patterns of gene expression via the
29 cooperative activities of DNA-bound TFs (Schmitz *et al.*, 2022). For proper activation of transcription,
30 the cell strictly regulates CRM activity by controlling TF access of CRM sequences through
31 nucleosome dynamics. Genome-wide approaches, such as Assay for Transposase Accessible
32 Chromatin sequencing (ATAC-seq), have been developed to profile accessible chromatin regions
33 (ACRs) (Buenrostro *et al.*, 2013; Minnoye *et al.*, 2021). In general, CRMs that localize to accessible
34 chromatin reflect active regulatory elements (Marand *et al.*, 2017; Schmitz *et al.*, 2022). Thus,
35 activation and silencing of gene expression is effectively controlled by the relative chromatin
36 accessibility of cognate CRMs.

37 CREs can be classified into distinct functional groups based on their regulatory effect on
38 transcription, including enhancers, silencers, promoters, and insulators (Schmitz *et al.*, 2022). Of
39 these, enhancers are of particular interest due to their transcription activating properties that function

40 independent of location and orientation of their target genes, in contrast to the stereotypical locations
41 of promoters surrounding gene transcription start sites (TSSs) (Marand *et al.*, 2017; Schmitz *et al.*,
42 2022). While analysis of chromatin accessibility in distinct tissues and cell types has been central to
43 identification of CRMs (Marand *et al.*, 2021), chromatin profiling techniques are largely qualitative and
44 lack the ability to quantitatively estimate regulatory activity. To overcome these challenges, Massively
45 Parallel Reporter Assays (MPRA) have been developed to quantify the transcription activating
46 properties of diverse sequences (Melnikov *et al.*, 2012; Arnold *et al.*, 2013;). In particular, Self-
47 Transcribing Active Regulatory Region with sequencing (STARR-seq) demonstrates the greatest
48 potential for broad application by eliminating the need for homogenous cell lines available only in
49 mammalian models typical of other MPRA methods (Arnold *et al.*, 2013; Ricci *et al.*, 2019; Sun *et al.*,
50 2019; Jores *et al.*, 2020). Although STARR-seq was originally designed to profile the entire genome for
51 regulatory activity, recent implementations have successfully utilized ATAC-seq libraries as input
52 (ATAC-STARR-seq), reducing the search space to potential regulatory regions and offsetting
53 sequencing costs and library complexity requirements (**Figure 1**). Despite its promise as a powerful
54 approach towards understanding *cis*-regulatory activity, computational analysis of ATAC-STARR-seq
55 data remains challenging, particularly due to a lack of dedicated software and computational pipelines.

56 Here, I present a computational pipeline for analysis of ATAC-STARR-seq data generated in *Zea*
57 *mays* L., cultivar B73 (Ricci *et al.*, 2019). After processing and evaluation of data quality, I demonstrate
58 how ATAC-STARR-seq data analysis allows for the interrogation of new biological questions. The
59 pipeline can be run entirely from the code below or through freely available bash, perl, and R scripts
60 hosted at https://github.com/Bio-protocol/Maize_ATAC_STARR_seq.

61



62

63 **Figure 1. Schematic of ATAC-STARR-seq**

64 ATAC-STARR-seq begins by first generating an ATAC-seq library. The ATAC fragments are then
65 cloned into a reporter assay and transformed into maize protoplasts. After X hours, transformed
66 protoplasts are then split into two pools, the first for sequencing the input fragments (ATAC-seq DNA),
67 and the second for purifying transcribed (mRNA) ATAC-seq fragments that facilitate their own
68 transcription from the reporter construct. Raw sequenced reads for the ATAC-seq input and mRNA
69 output are processed and aligned to the maize reference genome and compared to provide estimates
70 of cis-regulatory activity.

71

72 **Equipment**

73

74 This pipeline assumes that a user has knowledge of shell commands and is comfortable working
75 on a Linux-based operating system.

76

77 1. Computational Requirements

78 The following procedure can be run on any Linux-like system. However, this protocol and
79 publicly available code is written for executing commands via a high-performance computing
80 (HPC) cluster managed by a SLURM scheduler. However, the code presented here can be
81 readily converted to TORQUE or other HPC systems. The pipeline assumes a working Perl
82 interpreter version 5.30.0 or greater, and R version 3.6.2 or greater.

83

84 **Software**

85

86 The following analytical procedure makes use of several standard computational tools that are
87 assumed to be available in the user's shell environment.

88

89 **Software**

- 90 1. *BWA MEM* (Li and Durbin, 2009); v0.7.17; <http://bio-bwa.sourceforge.net/bwa.shtml>
- 91 2. *SAMtools* (Li *et al.*, 2009); v1.14; <http://www.htslib.org>
- 92 3. *BEDtools* (Quinlan and Hall, 2010); v2.27.1; <https://bedtools.readthedocs.io/en/latest/>
- 93 4. *SRA-toolkit* (Leinonen *et al.*, 2011); v2.11.1; <https://github.com/ncbi/sra-tools>
- 94 5. *fastp* (Chen *et al.*, 2018); v0.20.0; <https://github.com/OpenGene/fastp>
- 95 6. *pigz* v2.4; <https://zlib.net/pigz/>
- 96 7. *MACS2* (Liu, 2014); v2.2.7.1; <https://pypi.org/project/MACS2/>
- 97 8. *UCSC binaries* (Kent *et al.*, 2010); v1.04.0;
98 http://hgdownload.soe.ucsc.edu/admin/exe/linux.x86_64/
- 99 9. *tabix* (Li, 2011); v0.2.6; <http://www.htslib.org/doc/tabix.html>
- 100 10. *IGV* (Thorvaldsdottir *et al.*, 2013); v2.11.1; <https://software.broadinstitute.org/software/igv>
- 101 11. *MEME* (Grant *et al.*, 2011); v5.4.1; <https://meme-suite.org/meme/index.html>
- 102 12. *CrossMap* (Zhao *et al.*, 2014); v0.5.1; <http://crossmap.sourceforge.net>
- 103 13. *DeepTools* (Ramirez *et al.*, 2014); v3.5.1;
104 <https://deeptools.readthedocs.io/en/develop/index.html>

105

106 **Input data**

107 The starting input for this computational pipeline uses paired-end sequencing data from an ATAC-
108 STARR-seq experiment performed on maize protoplasts (Ricci *et al.*, 2019). The ATAC-STARR-
109 seq experiment consisted of a DNA input (ATAC-seq library) and a mRNA readout (self-transcribed
110 regulatory regions) to identify genomic regions exhibiting transcription-activating regulatory activity.

111

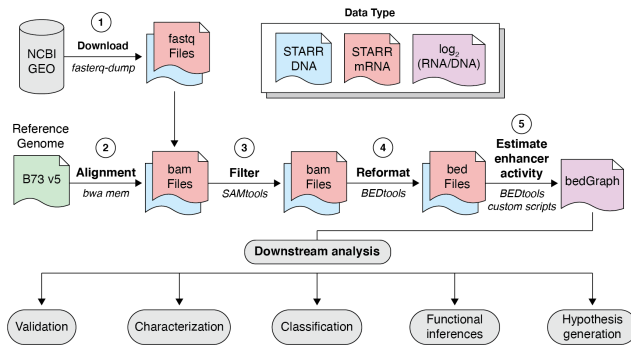
- 112 1. Transfected ATAC-seq DNA-input FASTQ
- 113 2. Transcribed ATAC-seq mRNA FASTQ

114

115 Procedure

116

117 An overview of the ATAC-STARR-seq pipeline is presented in **Figure 2**.



118

119 **Figure 2. Computational workflow for ATAC-STARR-seq data**

120 Raw sequence ATAC-STARR-seq data is first acquired from NCBI GEO, processed, and aligned to the
121 reference genome with *bwa mem*, filtered via *samtools*, reformatted as fragments with *bedtools*, and
122 compared via *bedtools* and custom *R* scripts to provide estimates of enhancer activity for downstream
123 analysis.

124

125 A. Download and prepare the requisite data and reference genome sequence

126 1. Raw mRNA ATAC-STARR-seq data generated from transfection of *Zea mays* leaf ATAC-seq
127 fragments in *Zea mays* protoplasts, and the accompanying ATAC-seq input fragments, are
128 publicly available on NCBI GEO. ATAC-STARR-seq mRNA and DNA input can be downloaded
129 with *fasterq-dump* available from the *SRA-Toolkit* package:

130

```
131 # set variables and download FASTQ files
```

```
132 mkdir FASTQ_files
```

```
133 cd FASTQ_files
```

```
134 fasterq-dump -o B73_maize_DNA_input.fastq SRR10964904
```

```
135 fasterq-dump -o B73_maize_mRNA_output.fastq SRR10964905
```

136

137 2. To save disk space, we will compress the FASTQ files with *pigz*. By default, *pigz* uses all
138 available processors or eight if the number of available processors is unknown. Alternatively,
139 users can use the unix tool, *gzip*, to compress the STARR mRNA and ATAC input DNA FASTQ
140 files.

141

```
142 # compress fastq files
```

```
143 pigz *.fastq
```

144

```
145 # NOT RUN
```

```
146 # Tip: gzip can be used as an alternative to pigz (parallel gzip)
```

147 # gzip *.fastq
148 3. Download the B73 reference genome sequence and gene annotation. The original article
149 mapped raw reads to version 4 of the B73 maize reference genome. In this case study, I map
150 and analyze maize ATAC-STARR-seq data to version 5 of the B73 maize reference genome to
151 showcase how updated reference genomes and read mapping strategies enable informative
152 reanalysis of publicly available data sets (Hufford *et al.*, 2021).

153
154 # download reference data
155 cd ../
156 mkdir Genome_Reference
157 cd Genome_Reference
158 wget [https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-](https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-NAM-5.0.fa.gz)
159 [NAM-5.0.fa.gz](https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-NAM-5.0.fa.gz)
160 wget [https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-](https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.gff3.gz)
161 [NAM-5.0_Zm00001eb.1.gff3.gz](https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-5.0/Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.gff3.gz)

162
163 4. To map data to the genome reference, we first need to decompress the FASTA file.
164 Constructing reference genome indices is a prerequisite for *BWA* alignment and allows for
165 faster post-alignment processing of BAM/SAM/BED formatted files with command line tools
166 such as *SAMtools* and *BEDtools*.

167
168 # create indices for reference genome FASTA
169 gunzip Zm-B73-REFERENCE-NAM-5.0.fa.gz
170 samtools faidx Zm-B73-REFERENCE-NAM-5.0.fa
171 bwa index Zm-B73-REFERENCE-NAM-5.0.fa

172 173 **B. Trim adapters and remove low quality reads**

174 1. Illumina platforms may produce reads with adapter sequences on the 3' ends if the DNA insert
175 is shorter than the number of cycles. Additionally, the fidelity of sequencing by synthesis
176 deteriorates with each additional cycle due to phasing, the desynchronization of cycles that
177 results from unremoved terminator caps ultimately leading to greater uncertainty of base calls
178 in later cycles. Removing adapter contamination and low-quality bases increases the total
179 number of alignable reads, particularly important when analyzing a relatively lower sequence
180 complexity experiment, such as ATAC-STARR-seq. We will use *fastp* to remove sequencing
181 adapters and low-quality reads for the mRNA output and DNA input. A script to perform read
182 trimming can be found here: [https://github.com/Bio-](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step01_trim_raw_reads.sh)
183 [protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step01_trim_raw_reads.sh](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step01_trim_raw_reads.sh)

184
185 # set variables

```
186     cd ../
187     fastqdir=$PWD/FASTQ_files
188     dna=B73_maize_DNA_input
189     rna=B73_maize_mRNA_output
190     threads=16
191
192     # trim and filter DNA input reads
193     fastp -j $dna.json -h $dna.html -w $threads \
194           -i $fastqdir/${dna}_1.fastq.gz -l $fastqdir/${dna}_2.fastq.gz \
195           -o $fastqdir/${dna}_1.trim.fastq.gz -O $fastqdir/${dna}_2.trim.fastq.gz
196
197     # trim and filter mRNA output reads
198     fastp -j $rna.json -h $rna.html -w $threads \
199           -i $fastqdir/${rna}_1.fastq.gz -l $fastqdir/${rna}_2.fastq.gz \
200           -o $fastqdir/${rna}_1.trim.fastq.gz -O $fastqdir/${rna}_2.trim.fastq.gz
201     sort -> $outdir/B73_maize_mRNA_output.raw.bam
202
203     # tidy up log files
204     mkdir fastp_log_files
205     mv *.json *.html fastp_log_files
206
```

207 C. Align and process sequenced reads

- 208 1. After trimming and quality filtering reads, we align the input DNA and output mRNA reads to
209 the maize B73 version 5 reference genome. To speed up the alignment and downstream
210 processing, we are using 24 CPUs (-t 24) to align the input and output reads. However, users
211 should modify this value to reflect the number of available cores on their system. Additionally,
212 we mark split hits as secondary alignments (-M) to be filtered out downstream as the maize
213 genome is highly repetitive. The output of *bwa mem* is piped to *samtools view* for compression
214 (SAM to BAM) and sorted by alignment coordinate prior to further processing to minimize the
215 footprint on disk. A script to perform these steps can be found here: https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step02_align_STARR_data.sh
216

```
217
218     cd ../
219     mkdir BAM_files
220     outdir=$PWD/BAM_files
221     refdir=$PWD/Genome_Reference
222     fastqdir=$PWD/FASTQ_files
223
224     # align DNA input and pipe to samtools for SAM to BAM conversion
```

```
225     bwa mem -M -t 24 $refdir/Zm-B73-REFERENCE-NAM-5.0.fa
226     $fastqdir/B73_maize_DNA_input_1.trim.fastq.gz
227     $fastqdir/B73_maize_DNA_input_2.trim.fastq.gz | samtools view -bS - | samtools sort - >
228     $outdir/B73_maize_DNA_input.raw.bam
229
230     # align RNA output and pipe to samtools for SAM to BAM conversion
231     bwa mem -M -t 24 $refdir/Zm-B73-REFERENCE-NAM-5.0.fa
232     $fastqdir/B73_maize_mRNA_output_1.trim.fastq.gz
233     $fastqdir/B73_maize_mRNA_output_2.trim.fastq.gz | samtools view -bS - | samtools sort - >
234     $outdir/B73_maize_mRNA_output.raw.bam
235
236     2. To ensure that only high quality alignments remain, here we remove non-properly paired reads
237     (-f 3), secondary hits (-F 256), alignments with low mapping quality (-q 10), and multiple
238     mapped reads (grep -v -E -e '\bXA:Z:') using a combination of samtools and unix commands.
239     The header, which contains information on the reference genome and the read mapping
240     parameters, is retained in the output by setting the -h flag in the samtools view command.
241
242     # filter DNA input alignments
243     samtools view -h -q 10 -f 3 -F 256 $outdir/B73_maize_DNA_input.raw.bam | grep -v -E -e
244     '\bXA:Z:' | samtools view -bSh - > $outdir/B73_maize_DNA_input.mq10.pp.unique.bam
245
246     # filter mRNA output alignments
247     samtools view -h -q 10 -f 3 -F 256 $outdir/B73_maize_mRNA_output.raw.bam | grep -v -E -e
248     '\bXA:Z:' | samtools view -bSh - > $outdir/B73_maize_mRNA_output.mq10.pp.unique.bam
249
250     3. A major difference between analysis of ATAC-seq and STARR-seq data is how assay
251     information is captured by sequencing. For paired-end ATAC-seq, since Tn5 inserts
252     sequencing adapters adjacent to its bound genomic location, chromatin accessibility is
253     encoded as the 5' ends of sequenced reads. In contrast, STARR-seq produces mRNA
254     transcripts from fragments that are capable of activating their own transcription, thus the entire
255     STARR-seq mRNA and DNA fragment is informative for analysis. The following commands
256     extract the coordinates of sequenced fragments by leveraging the CIGAR strings in BAM
257     paired-end alignments for DNA input and mRNA output. A script to extract fragments can be
258     found here: https://github.com/Bio-
259     protocol/Maize\_ATAC\_STARR\_seq/blob/master/workflow/step03\_extract\_fragments.sh
260
261     # create output directory
262     mkdir BED_files
263
```



```
264 # variables
265 outdir=$PWD/BAM_files
266 beddir=$PWD/BED_files
267 dna=B73_maize_DNA_input
268 rna=B73_maize_mRNA_output
269
270 # extract DNA input fragments (ignore the warnings from bedtools with respect to "missing"
271 mate pairs, these reflect one of the pairs that had its mate filtered in prior steps)
272 echo " extracting fragments from STARR DNA input ..."
273 samtools sort -n $outdir/$dna.mq10.pp.unique.bam \
274     | bedtools bamtobed -bedpe -i - \
275     | sort -k1,1 -k2,2n - \
276     | cut -f1,2,6 -> $beddir/$dna.fragments.bed
277
278 # extract mRNA output fragments
279 echo " extracting fragments from STARR mRNA output ..."
280 samtools sort -n $outdir/$rna.mq10.pp.unique.bam \
281     | bedtools bamtobed -bedpe -i - \
282     | sort -k1,1 -k2,2n - \
283     | cut -f1,2,6 -> $beddir/$rna.fragments.bed
284
```

285 **D. Identify regions with enriched activity over background**

```
286 1. To identify regions of the genome with the capacity to activate transcription, we assess
287 enrichment of mRNA reads relative to the input ATAC-seq fragments using macs2. As macs2
288 is a general peak caller, we need to adjust the default settings to tailor the analysis specifically
289 for ATAC-STARR-seq data. Since this experiment did not use unique molecular identifiers and
290 the number of transcripts from fragment is a direct reflection of its regulatory activity, duplicate
291 mRNA fragments are retained (--keep-dup all). To directly use coverages determined by the
292 input fragment coordinates, we turn off the default fragment shifting model (--nomodel) and set
293 the input type to BEDPE (-f BEDPE). Additionally, we reduce the maximum gap size between
294 candidate peaks to allow for the identification of fine-mapped regulatory elements within a
295 broader regulatory region (--max-gap 50) by setting the minimum peak size to 300 (--min-
296 length 300). Finally, we use the background coverage rates in place of the local bias which
297 aids in peak detection (--nolambda). A script to perform peak calling can be found here:
298 https://github.com/Bio-
299 protocol/Maize\_ATAC\_STARR\_seq/blob/master/workflow/step04\_call\_peaks.sh
```

```
300
301 # prepare output directory and input files
302 mkdir Peak_data
```

```
303
304     # generate input files
305     uniq $beddir/$rna.fragments.bed > $beddir/$rna.fragments.uniq.bed
306     uniq $beddir/$dna.fragments.bed > $beddir/$dna.fragments.uniq.bed
307     cat $beddir/$rna.fragments.uniq.bed $beddir/$dna.fragments.uniq.bed | sort -k1,1 -k2,2n - >
308     $beddir/ALL.fragments.uniq.bed
309
310     # find regulatory regions
311     echo " calling regulatory regions without duplicate removal ..."
312     macs2 callpeak -t $beddir/$rna.fragments.bed \
313         -c $beddir/ALL.fragments.uniq.bed \
314         --keep-dup all \
315         --max-gap 50 \
316         --min-length 300 \
317         --nolambda \
318         --nomodel \
319         -f BEDPE \
320         -g 1.6e9 \
321         --bdg \
322         -n STARR_wdups
323
324     # find regulatory regions
325     echo " calling regulatory regions with duplicate removal ..."
326     macs2 callpeak -t $beddir/$rna.fragments.uniq.bed \
327         -c $beddir/$dna.fragments.uniq.bed \
328         --keep-dup all \
329         --max-gap 50 \
330         --min-length 300 \
331         --nolambda \
332         --nomodel \
333         -f BEDPE \
334         -g 1.6e9 \
335         --bdg \
336         -n STARR_nodups
337
338     # find regulatory regions using all unique fragments
339     echo " calling regulatory regions by aggregating all unique fragments ..."
340     macs2 callpeak -t $beddir/ALL.fragments.uniq.bed \
341         --keep-dup all \
```

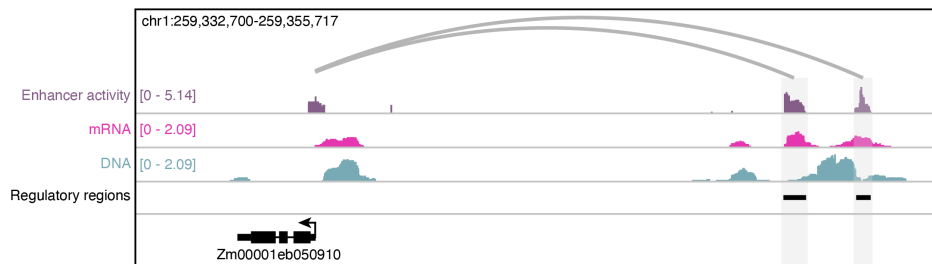
```
342         --max-gap 50 \  
343         --min-length 300 \  
344         --nolambda \  
345         --nomodel \  
346         -f BEDPE \  
347         -g 1.69e9 \  
348         --bdg \  
349         -n STARR_ALL  
350  
351     # clean-up output  
352     mv STARR_* Peak_data  
353  
354     # merge peaks  
355     cd Peak_data  
356     cat STARR_wdups_peaks.narrowPeak STARR_nodups_peaks.narrowPeak  
357     STARR_ALL_peaks.narrowPeak | sort -k1,1 -k2,2n - | bedtools merge -i - >  
358     STARR_merged_peaks.bed  
359  
360     2. To estimate regulatory activity at fine scale, first we need to create a list of unique intervals  
361     based on mRNA and DNA fragments. Next, we count the intersection of mRNA and DNA  
362     fragments for each unique interval. Finally, we remove all the temporary files to reduce the  
363     footprint on disk. A script to estimate enhancer activity can be found here:  
364     https://github.com/Bio-  
365     protocol/Maize\_ATAC\_STARR\_seq/blob/master/workflow/step05\_estimate\_enhancer\_activity.  
366     sh  
367  
368     # variables  
369     beddir=$PWD/BED_files  
370     dna=B73_maize_DNA_input  
371     ma=B73_maize_mRNA_output  
372     ref=./Genome_Reference/Zm-B73-REFERENCE-NAM-5.0.fa.fai  
373  
374     # sort the reference  
375     sort -k1,1 -k2,2n $ref > $ref.sorted  
376  
377     # merge RNA/DNA  
378     cat $beddir/$rna.fragments.uniq.bed $beddir/$dna.fragments.uniq.bed \  
379         | sort -k1,1 -k2,2n - \  
380         | bedtools genomecov -i - -bga -g $ref.sorted \  

```

```
381         | sort -k1,1 -k2,2n - \
382         | cut -f1-3 -> $beddir/Unique_genomic_intervals.bed
383
384     # count fragments
385     bedtools intersect -a $beddir/Unique_genomic_intervals.bed \
386         -b $beddir/$rna.fragments.bed \
387         -c -sorted -g $ref.sorted > $beddir/$rna.activity.raw.bed
388
389     bedtools intersect -a $beddir/$rna.activity.raw.bed \
390         -b $beddir/$dna.fragments.bed \
391         -c -sorted -g $ref.sorted > $beddir/B73_maize_mRNA_DNA.activity.raw.bed
392
393     # clean up temporary files
394     rm $beddir/Unique_genomic_intervals.bed
395     rm $beddir/$rna.activity.raw.bed
396
397     3. We and others define enhancer activity as the enrichment of mRNA transcripts that are
398         produced by DNA fragments in the assay in terms of  $\log_2(\text{mRNA}/\text{DNA})$  at unique fragment
399         intervals. Prior to taking the  $\log_2$  ratio of mRNA to DNA, we normalize both the input and
400         output to per million to account for differences in sequencing depth and complexity. A
401         pseudocount of one is added to any interval with at least one RNA or DNA fragment. The
402         following code can also be run from the command line using >Rscript
403         Estimate_Enhancer_Activity.R with the following script: https://github.com/Bio-
404         protocol/Maize\_ATAC\_STARR\_seq/blob/master/workflow/bin/Estimate\_Enhancer\_Activity.R
405
406         # open an interactive R session to estimate enhancer activity
407         cd $beddir
408         R
409
410         # load data
411         a <- read.table("B73_maize_mRNA_DNA.activity.raw.bed")
412
413         # reformat
414         rownames(a) <- paste(a$V1,a$V2,a$V3,sep="_")
415         a[,1:3] <- NULL
416         a <- as.matrix(a)
417         colnames(a) <- c("mRNA", "DNA")
418         a <- a[rowSums(a)!=0,]
419         a <- a + 1
```

```
420
421     # normalize
422     a <- a %*% diag(x=1e6/colSums(a))
423     colnames(a) <- c("mRNA", "DNA")
424     a <- as.data.frame(a)
425
426     # estimate enhancer activity
427     a$enhancer_activity <- log2(a$mRNA/a$DNA)
428
429     # reformat output
430     rownames(a) <- gsub("scaf_", "scaf", rownames(a))
431     df <- data.frame(do.call(rbind, strsplit(rownames(a), "_")))
432     df$X1 <- gsub("scaf", "scaf_", as.character(df$X1))
433     mrna <- df
434     dna <- df
435     df$X4 <- a$enhancer_activity
436     mrna$X4 <- a$mRNA
437     dna$X4 <- a$DNA
438
439     # cap negative activity at 0
440     df$X4 <- ifelse(df$X4 < 0, 0, df$X4)
441
442     # save enhancer activity BEDGRAPH file to disk
443     write.table(df, file="B73_maize.enhancer_activity.bdg", quote=F, row.names=F, col.names=F,
444     sep="\t")
445     write.table(mrna, file="B73_maize.mRNA.bdg", quote=F, row.names=F, col.names=F, sep="\t")
446     write.table(dna, file="B73_maize.DNA.bdg", quote=F, row.names=F, col.names=F, sep="\t")
447
448     # exit interactive mode
449     q()
450
451     4. To visualize enhancer activity and the normalized mRNA and DNA fragments at any given
452     locus, per million coverage values in bedGraph format from the previous step can be
453     converted into bigwig files (using bedGraphToBigWig from UCSC Utils) for facile visualization
454     using the Integrated Genomics Viewer (IGV) or JBrowse instances.
455
456     bedGraphToBigWig B73_maize.enhancer_activity.bdg ../Genome_Reference/Zm-B73-
457     REFERENCE-NAM-5.0.fai.sorted B73_maize.enhancer_activity.bw
458
```

- ```
459 bedGraphToBigWig B73_maize.mRNA.bdg ../Genome_Reference/Zm-B73-REFERENCE-
460 NAM-5.0.fa.fai.sorted B73_maize.mRNA.bw
461
462 bedGraphToBigWig B73_maize.DNA.bdg ../Genome_Reference/Zm-B73-REFERENCE-NAM-
463 5.0.fa.fai.sorted B73_maize.DNA.bw
464
465 5. To view the enhancer activity, mRNA, and DNA fragment bigwig files, download and install IGV
466 (https://software.broadinstitute.org/software/igv/download) on your local machine.
467
```



- 468 **Figure 3. Visualization of ATAC-STARR-seq data in *Zea mays***
- 469 Normalized (reads per million) coverages of the DNA ATAC-seq input (blue), self-transcribed
- 470 mRNA fragments (pink), and enhancer activity (purple;  $\log_2[\text{mRNA}/\text{DNA}]$ ) of a 23-kb window.
- 471 Regulatory regions are shown as black bars, while the grey loops reflect predicted enhancer-
- 472 gene links.
- 473
- 474
- 475 6. Unpack, bgzip, and index the gene product annotation. Then load all bigwig, narrowPeak, and
- 476 genome annotation files using “File > Load from File...” in IGV. An example of an IGV
- 477 screenshot is shown in **Figure 3**.

```
478
479 # change directory
480 cd ../Genome_Reference
481
482 # unzip
483 gunzip Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.gff3.gz
484
485 # sort by coordinate and remove whole chromosome intervals
486 sort -k1,1 -k4,4n Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.gff3 | grep -v assembly - >
487 Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.sorted.gff3
488
489 # compress with bgzip
490 bgzip Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.sorted.gff3
491
492 # index with tabix
```

493 tabix -p gff Zm-B73-REFERENCE-NAM-5.0\_Zm00001eb.1.sorted.gff3.gz

494

#### 495 **E. Create a list of control regions**

496 To assess enrichment of STARR regulatory regions determined by MACS2 relative to random control  
497 regions, we first need to identify regions of the genome that can be uniquely mapped given the  
498 sequencing output and read lengths. Although there are numerous methods for estimating mappability,  
499 I illustrate a simple approach using synthetic reads tailored to the sequencing parameters of the  
500 present experiment. First, we generate the same number of random read pairs with the same  
501 sequencing length (36 nucleotides) for mRNA and DNA input using the *wgsim* tool that is supplied to  
502 *SAMtools*. The synthetic reads are then remapped and uniformly processed as the original STARR-  
503 seq sequencing experiments to identify regions that are uniquely mappable. By constraining  
504 randomized control region selection to uniquely mappable genomic intervals, we ensure that  
505 downstream comparisons will not be biased by mappability and repeat composition. A script to  
506 construct control regions can be found here: [https://github.com/Bio-](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step06_create_control_regions.sh)  
507 [protocol/Maize\\_ATAC\\_STARR\\_seq/blob/master/workflow/step06\\_create\\_control\\_regions.sh](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step06_create_control_regions.sh)

508

509 # move into the "Genome\_Reference" directory

510 cd ./Genome\_Reference

511

512 # estimate read counts

513 mRNA\_counts=\$(samtools view -c ./BAM\_files/B73\_maize\_mRNA\_output.raw.bam)

514 DNA\_counts=\$(samtools view -c ./BAM\_files/B73\_maize\_DNA\_input.raw.bam)

515

516 # generate simulated reads matching the mRNA output using wgsim from the SAMtools  
517 package

518 wgsim -1 36 -2 36 -d 300 -N \$mRNA\_counts Zm-B73-REFERENCE-NAM-5.0.fa

519 simulated\_STARR\_mRNA\_r1.fq simulated\_STARR\_mRNA\_r2.fq

520

521 # generate simulated reads matching the DNA input

522 wgsim -1 36 -2 36 -d 300 -N \$DNA\_counts Zm-B73-REFERENCE-NAM-5.0.fa

523 simulated\_STARR\_DNA\_r1.fq simulated\_STARR\_DNA\_r2.fq

524

525 # compress synthetic fastq files

526 pigz \*.fq

527

528 1. Remap the synthetic reads using the same pipeline as for the original STARR-seq data.

529

530 # variables

531 outdir=\$PWD/BAM\_files

```
532 refdir=$PWD/Genome_Reference
533 ref=$refdir/Zm-B73-REFERENCE-NAM-5.0.fa.fai.sorted
534 fastqdir=$refdir
535
536 # align synthetic mRNA output and pipe to samtools for SAM to BAM conversion
537 bwa mem -M -t 24 $refdir/Zm-B73-REFERENCE-NAM-5.0.fa \
538 $fastqdir/simulated_STARR_mRNA_r1.fq.gz \
539 $fastqdir/simualted_STARR_mRNA_r2.fq.gz \
540 | samtools view -bS - \
541 | samtools sort - > $outdir/simulated_STARR_mRNA.raw.bam
542
543 # align synthetic DNA input and pipe to samtools for SAM to BAM conversion
544 bwa mem -M -t 24 $refdir/Zm-B73-REFERENCE-NAM-5.0.fa \
545 $fastqdir/simulated_STARR_DNA_r1.fq.gz \
546 $fastqdir/simulated_STARR_DNA_r2.fq.gz \
547 | samtools view -bS - \
548 | samtools sort - > $outdir/simulated_STARR_DNA.raw.bam
549
550 2. Process and filter reads using the original STARR-seq pipeline.
551
552 # filter synthetic mRNA alignments
553 echo " filtering synthetic STARR mRNA alignments ..."
554 samtools view -h -q 10 -f 3 $outdir/simulated_STARR_mRNA.raw.bam \
555 | grep -v -E -e '\bXA:Z:' \
556 | samtools view -bSh - > $outdir/simulated_STARR_mRNA.mq10.pp.unique.bam
557
558 # filter synthetic DNA alignments
559 echo " filtering STARR DNA alignments ..."
560 samtools view -h -q 10 -f 3 $outdir/simulated_STARR_DNA.raw.bam \
561 | grep -v -E -e '\bXA:Z:' \
562 | samtools view -bSh - > $outdir/simulated_STARR_DNA.mq10.pp.unique.bam
563
564 3. Identify uniquely mappable regions.
565
566 # extract mRNA output fragments
567 echo " extracting fragments from simulated STARR mRNA output ..."
568 samtools sort -n $outdir/simulated_STARR_mRNA.mq10.pp.unique.bam \
569 | bedtools bamtoBED -bedpe -i - \
570 | sort -k1,1 -k2,2n - \
```

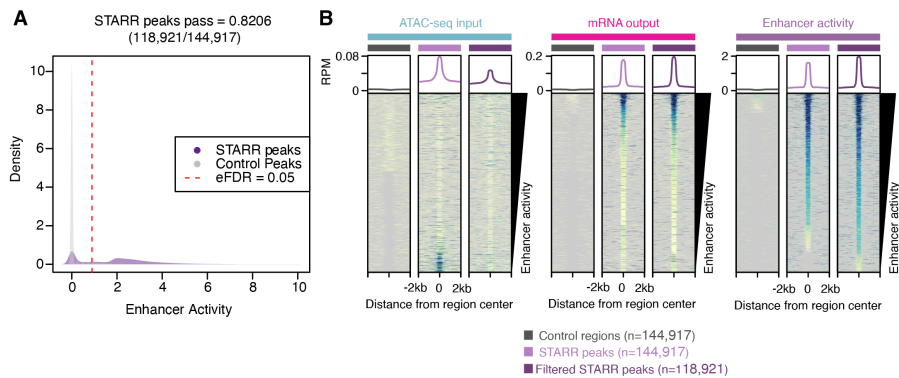


```
571 | cut -f1,2,6 -> $refdir/simulated_STARR_mRNA.fragments.bed
572
573 # extract DNA input fragments
574 echo " extracting fragments from simulated STARR DNA input ..."
575 samtools sort -n $outdir/simulated_STARR_DNA.mq10.pp.unique.bam \
576 | bedtools bamtobed -bedpe -i - \
577 | sort -k1,1 -k2,2n - \
578 | cut -f1,2,6 -> $refdir/simulated_STARR_DNA.fragments.bed
579
580 # merge all fragments (sorting by coordinate at this step may take a while)
581 cat $refdir/simulated_STARR_mRNA.fragments.bed
582 $refdir/simulated_STARR_DNA.fragments.bed \
583 | sort -k1,1 -k2,2n \
584 | bedtools merge -i -> $refdir/mappable_genomic_regions.bed
585
586 4. Construct control regions using only mappable regions and excluding putative regulatory
587 regions output by MACS2.
588
589 # create controls
590 peaks=$PWD/Peak_data/STARR_merged_peaks.bed
591 bedtools shuffle -i $peaks \
592 -g $ref \
593 -incl $refdir/mappable_genomic_regions.bed \
594 -excl $peaks \
595 | sort -k1,1 -k2,2n -> $PWD/Peak_data/STARR_CONTROL.bed
596
597 F. Compare enhancer activity
598 1. Determine enhancer activity for predicted enhancers output by MACS2 as well as the negative
599 control regions.
600
601 # create directory to contain analysis
602 cd ../
603 mkdir 01_Peak_Analysis
604 cd 01_Peak_Analysis
605
606 # map maximum enhancer activity to putative regulatory regions (wdups)
607 bedtools map -a ../Peak_data/STARR_merged_peaks.bed -
608 b ../BED_files/B73_maize.enhancer_activity.bdg -o max -c 4 >
609 STARR_merged_peaks.enhancer_activity.bed
```

```

610
611 # map maximum enhancer activity to control
612 bedtools map -a ../Peak_data/STARR_CONTROL.bed -
613 b ../BED_files/B73_maize.enhancer_activity.bdg -o max -c 4 >
614 STARR_CONTROL.enhancer_activity.bed
615

```



616  
617 **Figure 4. Analysis of STARR regulatory region enhancer activity.**

618 (A) Distribution of enhancer activity for STARR peaks (purple) and random control regions  
619 (grey). Dashed red line indicates the 95% quantile of enhancer activity of random control  
620 regions. (B) Average (top) and individual site heatmaps of reads per million (RPM) for ATAC-  
621 seq input (left), mRNA output (middle) and enhancer activity (right) for control regions, all  
622 STARR peaks, and the filtered STARR peak set.

623  
624 2. To remove regulatory regions with enhancer activity similar to background, we filter STARR  
625 regulatory regions using an empirical false discovery rate (eFDR) based on the matched  
626 control regions. A user-specified eFDR threshold identifies the enhancer activity value in the  
627 control set that removes 1-eFDR percent of control regions. In this example, we set the FDR  
628 to a widely used rate of 0.05. The following code performs and plots eFDR filtering and  
629 enhancer activity distributions and can be run from the command line using *>Rscript*  
630 *eFDR\_Filter\_STARR\_Peaks.R*. Filtering STARR peaks based on eFDR thresholds derived  
631 from the control regions is visualized in **Figure 4A**. An R script of the following code can be  
632 found here: [https://github.com/Bio-](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/bin/eFDR_Filter_STARR_Peaks.R)  
633 [protocol/Maize\\_ATAC\\_STARR\\_seq/blob/master/workflow/bin/eFDR\\_Filter\\_STARR\\_Peaks.R](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/bin/eFDR_Filter_STARR_Peaks.R)

```

634
635 # start an interactive R session
636 > R
637
638 # load libraries
639 library(scales)
640
641 # load data

```

```
642 starr <- read.table("STARR_merged_peaks.enhancer_activity.bed")
643 con <- read.table("STARR_CONTROL.enhancer_activity.bed")
644
645 # set missing to 0
646 starr$V4[starr$V4=='.'] <- 0
647 con$V4[con$V4=='.'] <- 0
648
649 # convert to numeric
650 starr$V4 <- as.numeric(as.character(starr$V4))
651 con$V4 <- as.numeric(as.character(con$V4))
652
653 # get empirical thresholds
654 fdr <- 0.05
655 threshold <- quantile(con$V4, (1-fdr))
656
657 # filter STARR regulatory regions
658 filtered <- subset(starr, starr$V4 >= threshold)
659
660 # estimate fraction of retained regions
661 frac <- signif(nrow(filtered)/nrow(starr), digits=4)
662
663 # set up multipanel plot area
664 pdf("Density_eFDR_STARR_Peak_Filtering.pdf", width=5, height=5)
665
666 # plot control/observed enhancer activities for STARR peaks with duplicates
667 den.starr <- density(starr$V4)
668 den.con <- density(con$V4)
669 plot(NA,
670 xlab="Enhancer Activity",
671 ylab="Density",
672 xlim=c(range(range(den.starr$x), range(den.con $x))),
673 ylim=c(range(range(den.starr$y), range(den.con$y))))
674 polygon(x=c(min(den.starr$x), den.starr$x, max(den.starr$x)),
675 y=c(0, den.starr$y, 0), col=alpha("darkorchid4", 0.5), border=NA)
676 polygon(x=c(min(den.con$x), den.con$x, max(den.con$x)),
677 y=c(0, den.con$y, 0), col=alpha("grey80", 0.5), border=NA)
678 abline(v=threshold, col="red", lty=2)
679 mtext(paste0("STARR peaks pass = ",frac," (", nrow(filtered), "/", nrow(starr),")"))
680
```

```
681 legend("right", legend=c("STARR Peaks", "Control Peaks", paste0("eFDR = ", fdr)),
682 col=c("darkorchid4", "grey75", "red"), border=c(NA, NA, "red"), pch=c(16, 16, NA), lty=c(NA,
683 NA, 2))
684
685 # close device
686 dev.off()
687
688 # save filtered STARR regulatory regions
689 write.table(filtered, file="STARR_merged_peaks.enhancer_activity.eFDR05.bed", quote=F,
690 row.names=F, col.names=F, sep="\t")
691
692 # exit R
693 q()
694
695 3. We can now assess the relative enhancer activities across all regions for the filtered and
696 unfiltered STARR peaks and controls using DeepTools (Figure 4B). A script to plot heatmaps
697 via DeepTools can be found here: https://github.com/Bio-
698 protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step07_plot_enhancer_activity.sh
699
700 # load function
701 getmaps(){
702
703 # input
704 ina=$1
705 id=$2
706 dat=../BED_files/*.bw
707
708 # output
709 outa=$id.mat.gz
710 outm=$id.mat.txt
711 fig=$id.pdf
712
713 # parameters
714 threads=48
715 window=2000
716 bin=20
717 cols=YIGnBu
718
719 # create matrix
```

```
720 computeMatrix reference-point --referencePoint center \
721 -S $dat \
722 -b $window -a $window \
723 -R $ina \
724 --missingDataAsZero \
725 -o $outa \
726 --outFileNameMatrix $outm \
727 -p $threads --binSize $bin
728
729 # plot heatmap
730 plotHeatmap --matrixFile $outa \
731 --colorMap YlGnBu \
732 -out $id.heatmap.pdf
733
734 }
735 export -f getmaps
736
737 # run for each file
738 getmaps STARR_merged_peaks.enhancer_activity.bed STARR_peaks
739 getmaps STARR_merged_peaks.enhancer_activity.eFDR05.bed STARR_peaks_filtered
740 getmaps STARR_CONTROL.enhancer_activity.bed control_regions
741
```

## 742 **G. Identification of large regulatory domains in the maize genome**

743 1. One question these data allow us to ask is whether a relationship exists between the size of a  
744 regulatory region and its enhancer activity. So called “super enhancers” in mammalian  
745 systems describe hyperactive transcription-activating regulatory domains associated with cell  
746 identity that exhibit increased density of TF binding sites compared to typical enhancers  
747 (Hnisz *et al.*, 2013). Integration of the STARR peaks and enhancer activities with other data  
748 sets allows us to determine whether similar hyperactive regulatory domains exist in maize. To  
749 query TF binding site density, we first download position weight matrices of known TFs from  
750 the *meme* database and identify putative TF binding sites using *fimo* (also from the *meme*  
751 suite) conditioning on a *P*-value threshold less than 1e-5. A script to identify large regulatory  
752 domains can be found here: [https://github.com/Bio-](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step08_identify_large_regulatory_domains.sh)  
753 [protocol/Maize\\_ATAC\\_STARR\\_seq/blob/master/workflow/step08\\_identify\\_large\\_regulatory\\_do](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step08_identify_large_regulatory_domains.sh)  
754 [mains.sh](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq/blob/master/workflow/step08_identify_large_regulatory_domains.sh)  
755  
756 # make a new directory for the TFBS analysis  
757 cd ../  
758 mkdir 02\_Hyperactive\_Regulatory\_Region\_Analysis

```
759 cd ./02_Hyperactive_Regulatory_Region_Analysis
760
761 # download and decompress motif databases
762 wget https://meme-suite.org/meme/meme-
763 software/Databases/motifs/motif_databases.12.23.tgz
764 tar -xvzf motif_databases.12.23.tgz
765 rm motif_databases.12.23.tgz
766
767 # variables
768 threads=16
769 ref=./Genome_Reference/Zm-B73-REFERENCE-NAM-5.0.fa
770 peaks=./01_Peak_Analysis/STARR_merged_peaks.enhancer_activity.eFDR05.bed
771 controls=./01_Peak_Analysis/STARR_CONTROL.enhancer_activity.bed
772 motifs=./motif_databases/ARABD/ArabidopsisDAPv1.meme
773
774 # extract fasta sequences
775 bedtools getfasta -bed $peaks -fi $ref -fo $peaks.fasta
776 bedtools getfasta -bed $controls -fi $ref -fo $controls.fasta
777
778 # identify putative TFBS
779 fimo --oc TFBS_peaks $motifs $peaks.fasta
780 fimo --oc TFBS_controls $motifs $controls.fasta
781
782 # reformat fimo output (filtering p-value > 1e-5) using the perl script provided in the github
783 repository (https://github.com/Bio-
784 protocol/Maize_ATAC_STARR_seq/blob/master/workflow/bin/convertMotifCoord.pl)
785 perl convertMotifCoord.pl TFBS_peaks/fimo.gff | sed -e 's/_tnt/g' - | sort -k1,1 -k2,2n - >
786 TFBS_peaks.motifs.bed
787 perl convertMotifCoord.pl TFBS_controls/fimo.gff | sed -e 's/_tnt/g' - | sort -k1,1 -k2,2n - >
788 TFBS_controls.motifs.bed
789
790 # annotate motif coverage/counts for STARR and control peaks
791 bedtools annotate -
792 i ./01_Peak_Analysis/STARR_merged_peaks.enhancer_activity.eFDR05.bed -files
793 TFBS_peaks.motifs.bed -both | sort -k1,1 -k2,2n - >
794 STARR_merged_peaks.enhancer_activity.eFDR05.ann.bed
795 bedtools annotate -i ./01_Peak_Analysis/STARR_CONTROL.enhancer_activity.bed -files
796 TFBS_controls.motifs.bed -both | sort -k1,1 -k2,2n - >
797 STARR_CONTROL.enhancer_activity.ann.bed
```

```
798
799 # extract genes
800 perl -ne 'if($_ =~ /^#/){next;}chomp;my@col=split("\t",$_);if($col[2] eq
801 'gene'){print"$_\n";}' ../Genome_Reference/Zm-B73-REFERENCE-NAM-
802 5.0_Zm00001eb.1.gff3 | sort -k1,1 -k4,4n - > ../Genome_Reference-
803 NAM-5.0_Zm00001eb.1.genes.gff3
804
805 # classify genomic context of STARR and control peaks (you can ignore the warnings from
806 bedtools about inconsistent naming conventions, you can thank the genome assembly team
807 for these annoying, but harmless warnings)
808 bedtools closest -a STARR_merged_peaks.enhancer_activity.eFDR05.ann.bed -
809 b ../Genome_Reference/Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.genes.gff3 -D b >
810 STARR_merged_peaks.enhancer_activity.eFDR05.ann2.bed
811 bedtools closest -a STARR_CONTROL.enhancer_activity.ann.bed -
812 b ../Genome_Reference/Zm-B73-REFERENCE-NAM-5.0_Zm00001eb.1.genes.gff3 -D b >
813 STARR_CONTROL.enhancer_activity.ann2.bed
814
815 # clean up
816 mv STARR_merged_peaks.enhancer_activity.eFDR05.ann2.bed
817 STARR_merged_peaks.enhancer_activity.eFDR05.ann.bed
818 mv STARR_CONTROL.enhancer_activity.ann2.bed
819 STARR_CONTROL.enhancer_activity.ann.bed
820
821 2. We can now investigate the relationship among regulatory region size, motif density, and
822 enhancer activity to identify putative regulatory domains (Figure 5A-5F). To do so, we will start
823 an interactive R session and load the annotated peak and control files from above. A script to
824 automate the following code can be found here: https://github.com/Bio-
825 protocol/Maize_ATAC_STARR_seq/blob/master/workflow/bin/Characterize_Regulatory_Regio
826 ns.R
827
828 # open R
829 > R
830
831 # Analyze regulatory regions
832
833 # load libraries
834 library(vioplots)
835 library(dplyr)
836 library(MASS)
```

```
837 library(RColorBrewer)
838 library(scales)
839
840 # load data
841 starr <- read.table("STARR_merged_peaks.enhancer_activity.eFDR05.ann.bed")
842 control <- read.table("STARR_CONTROL.enhancer_activity.ann.bed")
843
844 # select random control regions to match the filtered STARR peaks
845 control <- control[sample(nrow(starr)),]
846
847 # rename columns for clarity (frac_RR_motif = fraction of regulatory region covered by motifs)
848 starr[,7:15] <- NULL
849 control[,7:15] <- NULL
850 colnames(starr)[4:7] <- c("activity", "motif_counts", "frac_RR_motif", "gene_distance")
851 colnames(control)[4:7] <- c("activity", "motif_counts", "frac_RR_motif", "gene_distance")
852
853 # classify
854 starr$class <- ifelse((starr$gene_distance < 0 & starr$gene_distance > -200), "TSS",
855 ifelse(starr$gene_distance < -200 & starr$gene_distance > -2000, "promoter",
856 ifelse(starr$gene_distance > 0 & starr$gene_distance < 1000, "TTS",
857 ifelse(starr$gene_distance == 0, "genic", "intergenic"))))
858
859 control$class <- ifelse((control$gene_distance < 0 & control$gene_distance > -200), "TSS",
860 ifelse(control$gene_distance < -200 & control$gene_distance > -2000, "promoter",
861 ifelse(control$gene_distance > 0 & control$gene_distance < 1000, "TTS",
862 ifelse(control$gene_distance == 0, "genic", "intergenic"))))
863
864 # plot distribution
865 pdf("STARR_peak_control_genomic_distribution.pdf", width=10, height=5)
866 layout(matrix(c(1:2), nrow=1))
867 pie(table(starr$class))
868 pie(table(control$class))
869 dev.off()
870
871 # estimate regulatory region size (in log10 scale)
872 starr$size <- log10(starr$V3-starr$V2)
873 control$size <- log10(control$V3-control$V2)
874
875 # compare sizes between peaks and controls (sanity check)
```

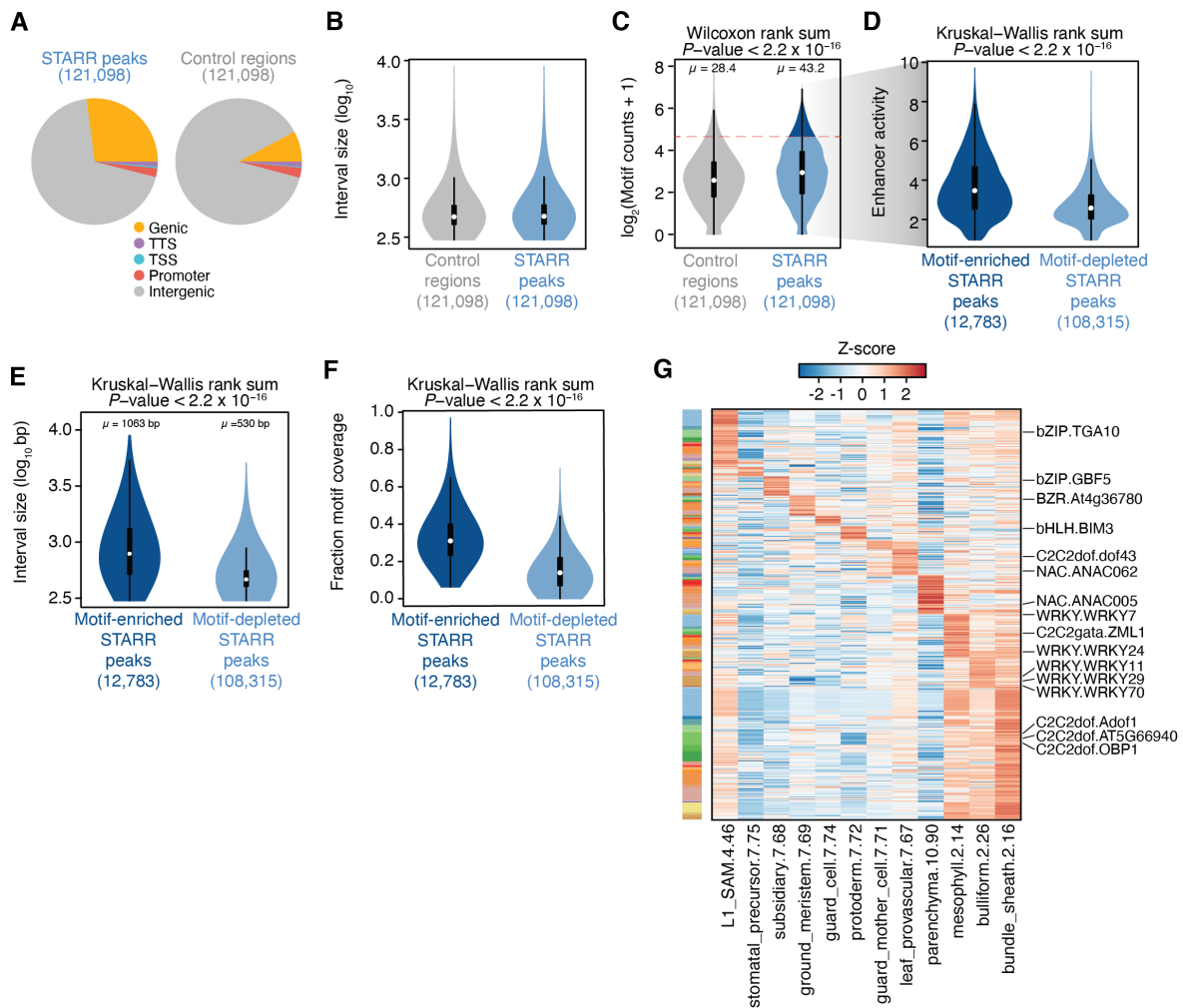


```
876 pdf("STARR_peak_control_sizes.pdf", width=5, height=6)
877 vioplot(starr$size, control$size,
878 ylab="Interval size (log10)",
879 col=c("dodgerblue", "grey75"),
880 names=c(paste0("STARR peaks \n (n=",nrow(starr),"),",
881 paste0("Control regions \n (n=",nrow(control),")"))))
882 dev.off()
883
884 # compare motif counts
885 pval <- wilcox.test(starr$motif_counts, control$motif_counts)$p.value
886 pval <- ifelse(pval==0, 2.2e-16, pval)
887 mean.peak <- mean(starr$motif_counts)
888 mean.cont <- mean(control$motif_counts)
889
890 # find 95% quantile for control motif count
891 upper.threshold <- quantile(control$motif_counts, 0.95)
892
893 # plot
894 pdf("STARR_peak_control_motif_counts.pdf", width=5, height=6)
895 vioplot(log1p(starr$motif_counts), log1p(control$motif_counts),
896 ylab="log2(Motif counts + 1)",
897 col=c("dodgerblue", "grey75"),
898 names=c(paste0("STARR peaks \n (n=",nrow(starr),"),",
899 paste0("Control regions \n (n=",nrow(control),")"))),
900 ylim=c(0,8),
901 areaEqual=T,
902 h=0.25)
903 mtext(paste0("Wilcoxon Rank Sum P-value = ", signif(pval, digits=3)))
904 text(1, 7.5, labels=paste0("Mean = ", signif(mean.peak, digits=3)))
905 text(2, 7.5, labels=paste0("Mean = ", signif(mean.cont, digits=3)))
906 points(1, log1p(upper.threshold), col="red", pch="-")
907 points(2, log1p(upper.threshold), col="red", pch="-")
908 dev.off()
909
910 # split STARR regions by motif counts based on 95% quantile control dist
911 starr$group <- ifelse(starr$motif_counts >= upper.threshold, "high", "low")
912 pval <- kruskal.test(starr$activity, starr$group)$p.value
913 pdf("STARR_peak_activity_vs_group.pdf", width=5, height=6)
914 vioplot(starr$activity~starr$group,
```

```
915 ylab="Enhancer activity",
916 col=c("dodgerblue4", "dodgerblue"),
917 names=c(paste0("Motif-enriched \n STARR peaks \n (n=",
918 nrow(starr[starr$group=="high",]),"),"),
919 paste0("Motif-depleted \n STARR peaks \n (n=",
920 nrow(starr[starr$group=="low",]),"),"),
921 areaEqual=F,
922 xlab="",
923 h=0.25)
924 mtext(paste0("Kruskal-Wallis rank sum P-value = ", signif(pval, digits=3)))
925 dev.off()
926
927 # compare STARR region size
928 pval <- kruskal.test(starr$size, starr$group)$p.value
929 pval <- ifelse(pval==0, 2.2e-16, pval)
930 pdf("STARR_peak_size_vs_group.pdf", width=5, height=6)
931 vioplot(starr$size~starr$group,
932 ylab="Interval size (log10)",
933 col=c("dodgerblue4", "dodgerblue"),
934 names=c(paste0("Motif-enriched \n STARR peaks \n (n=",
935 nrow(starr[starr$group=="high",]),"),"),
936 paste0("Motif-depleted \n STARR peaks \n (n=",
937 nrow(starr[starr$group=="low",]),"),"),
938 areaEqual=F,
939 xlab="",
940 h=0.25)
941 mtext(paste0("Kruskal-Wallis rank sum P-value = ", signif(pval, digits=3)))
942 dev.off()
943
944 # compare motif coverage
945 pval <- kruskal.test(starr$frac_RR_motif, starr$group)$p.value
946 pval <- ifelse(pval==0, 2.2e-16, pval)
947 pdf("STARR_peak_motif_coverage_vs_group.pdf", width=5, height=6)
948 vioplot(starr$frac_RR_motif~starr$group,
949 ylab="Fraction motif coverage",
950 col=c("dodgerblue4", "dodgerblue"),
951 names=c(paste0("Motif-enriched \n STARR peaks \n (n=",
952 nrow(starr[starr$group=="high",]),"),"),
953 paste0("Motif-depleted \n STARR peaks \n (n=",
```

```
954 nrow(starr[starr$group=="low",]),"")),
955 areaEqual=F,
956 xlab="")
957 mtext(paste0("Kruskal-Wallis rank sum P-value = ", signif(pval, digits=3)))
958 dev.off()
959
960 # split by group
961 starr.me <- subset(starr, starr$group=="high")
962 starr.md <- subset(starr, starr$group=="low")
963 write.table(starr.me,
964 file="STARR_starrs_peaks.enhancer_activity.eFDR05.ann.high_motif.bed",
965 quote=F, row.names=F, col.names=F, sep="\t")
966 write.table(starr.md,
967 file="STARR_starrs_peaks.enhancer_activity.eFDR05.ann.low_motif.bed",
968 quote=F, row.names=F, col.names=F, sep="\t")
969
970
```

971



**Figure 5: Identification of motif-dense enhancer regulatory domains.**

(A) Genomic distribution of STARR peaks (left) and control regions (right). (B) Distribution of control region (grey) and STARR peak (blue) interval lengths. (C) Distribution of motif counts in control regions (grey) and STARR peaks (blue). The dashed red line indicates the 95% quantile of motif counts from control regions used to classify STARR peaks into high and low motif count classes. (D) Distribution of enhancer activity for STARR peaks with enriched (dark blue) and depleted (light blue) motif counts. (E) Distribution of interval lengths for motif-enriched (dark blue) and motif-depleted (light blue) STARR peaks. (F) Distribution of fraction of STARR peak covered by motif for motif-enriched (dark blue) and motif-depleted (light blue) STARR peaks. (G) Heatmap illustrating Z-score transformed motif enhancer activities across intergenic motif-enriched STARR peaks scaled by the relative chromatin accessibility in various maize cell types.

972

973

974

- To determine if the large intergenic regulatory domain regions are associated with cell identity, we will compare enhancer activities versus various cell-type-specific accessible chromatin regions (ACRs) leveraging a recent single-cell ATAC-seq (scATAC-seq) dataset from multiple maize organs (Marand *et al.*, 2021). First, download the matrix containing normalized

975

976

977

```
978 accessibility counts across accessible chromatin regions for each profiled cell type. We then
979 extract ACR genomic coordinates (which are in version 4 of the B73 reference genome) and
980 convert them to version 5 of the B73 reference genome using the CrossMap tool and chain
981 file.
982
983 # download the counts matrix
984 wget -O maize_scATAC_atlas_ACR_celltype_CPM.txt.gz
985 https://www.ncbi.nlm.nih.gov/geo/download/?acc=GSE155178&format=file&file=GSE1551
986 78%5Fmaize%5FscATAC%5Fatlas%5FACR%5Fcelltype%5FCPM%2Etxt%2Egz
987
988 # unzip
989 gunzip maize_scATAC_atlas_ACR_celltype_CPM.txt.gz
990
991 # download chain file
992 wget https://download.maizegdb.org/Zm-B73-REFERENCE-NAM-
993 5.0/chain_files/B73_RefGen_v4_to_Zm-B73-REFERENCE-NAM-5.0.chain
994
995 # extract coordinates and conform chromosome names to V4 reference
996 cut -f1 maize_scATAC_atlas_ACR_celltype_CPM.txt \
997 | grep '^chr' - \
998 | perl -ne 'chomp;my@col=split("_",$_);print"$col[0]\t$col[1]\t$col[2]\n";' - \
999 | sed -e 's/chrB73V4ctg/B73V4_ctg/g' - \
1000 | sed -e 's/chr//g' - \
1001 | sort -k1,1 -k2,2n - > maize_scATAC_atlas_ACRs.bed
1002
1003 # convert ACR coordinates from V4 to V5
1004 CrossMap.py bed B73_RefGen_v4_to_Zm-B73-REFERENCE-NAM-5.0.chain
1005 maize_scATAC_atlas_ACRs.bed > maize_scATAC_atlas_ACRs.V4_to_V5.txt
1006
1007 # discard unmapped and split projections
1008 grep -v 'Unmap|split' maize_scATAC_atlas_ACRs.V4_to_V5.txt \
1009 | perl -ne 'chomp; my@col=split("\t",$_);
1010 print"chr$col[0]_$col[1]_$col[2]\tchr$col[4]_$col[5]_$col[6]\n";' - \
1011 | sort -k1,1 -k2,2n - \
1012 | sed -e 's/chrscaf/scaf/g' - \
1013 | sed -e 's/chrB73V4_ctg/chrB73V4ctg/g' - > maize_scATAC_atlas_ACRs.V4_to_V5.clean.txt
1014
1015 # update ACR coordinates in matrix file using R
1016 > R
```

```
1017
1018 # read into data frames
1019 conv <- read.table("maize_scATAC_atlas_ACRs.V4_to_V5.clean.txt")
1020 mat <- read.table("maize_scATAC_atlas_ACR_celltype_CPM.txt")
1021
1022 # subset mat rows by retained ACRs after projection
1023 shared <- intersect(rownames(mat), as.character(conv$V1))
1024 mat <- mat[shared,]
1025 rownames(conv) <- conv$V1
1026 conv <- conv[shared,]
1027
1028 # update mat rowIDs
1029 rownames(mat) <- conv$V2
1030
1031 # save output
1032 write.table(mat, file="maize_scATAC_atlas_ACR_celltype_CPM.V5.txt", quote=F,
1033 row.names=T, col.names=T, sep="\t")
1034
1035 # exit R
1036 q()
1037
1038 # remove temporary files
1039 rm maize_scATAC_atlas_ACR_celltype_CPM.txt maize_scATAC_atlas_ACRs.bed
1040 maize_scATAC_atlas_ACRs.V4_to_V5.txt maize_scATAC_atlas_ACRs.V4_to_V5.clean.txt
1041
1042
1043 4. Intersect the scATAC-seq ACRs with the STARR peaks with enriched motif counts. Load the
1044 scATAC-seq matrix and intersected ACRs/STARR peaks files into R to estimate enhancer
1045 activity enrichment scaled by relative accessibilities across cell types. As the STARR-seq data
1046 was derived from maize seedlings, we will further restrict the analysis of scATAC-seq cell types
1047 to those derived primarily from maize seedlings. The following code written in R can be
1048 executed with the script named 'motif_enhancer_activity_maize_celltypes.R' and provides
1049 estimates of enhancer activity over various motifs scaled by the relative cell type accessibility,
1050 allowing insights into cell-type-specific transcription factor regulation of active enhancers
1051 (Figure 5G).
1052
1053 # extract ACR coordinates
1054 cut -f1 maize_scATAC_atlas_ACR_celltype_CPM.V5.txt | grep -v 'unknown.5.50' | sed -e
1055 's/scaf_/scaf/g' | perl -ne 'chomp;my@col=split(" ",$_);print"$col[0]\t$col[1]\t$col[2]\n";' - | sed -
```

```
1056 e 's/scaf/scaf_/g' - | sort -k1,1 -k2,2n - > maize_scATAC_atlas_ACRs.V5.bed
1057
1058 # intersect scATAC ACRs with high motif counts STARR peaks
1059 bedtools intersect -a STARR_starrs_peaks.enhancer_activity.eFDR05.ann.high_motif.bed -b
1060 maize_scATAC_atlas_ACRs.V5.bed -wa -wb >
1061 STARR_starrs_peaks.enhancer_activity.eFDR05.ann.high_motif.scATAC_ACRs.bed
1062
1063 # map enhancer activity over motifs
1064 bedtools map -a TFBS_peaks.motifs.bed -b ../BED_files/B73_maize.enhancer_activity.bdg -c
1065 4 -o max > TFBS_peaks.motifs.enhancer_activity.bed
1066
1067 # motifs to large regulatory regions
1068 bedtools intersect -a TFBS_peaks.motifs.enhancer_activity.bed -b
1069 STARR_starrs_peaks.enhancer_activity.eFDR05.ann.high_motif.scATAC_ACRs.bed -wa -wb
1070 > TFBS_peaks.motifs.enhancer_activity.bed
1071
1072 # open R (alternatively, a script to automate the following code can be found here:
1073 https://github.com/Bio-
1074 protocol/Maize_ATAC_STARR_seq/blob/master/workflow/bin/motif_enhancer_activity_maize_
1075 celltypes.R)
1076 > R
1077
1078 # estimate enhancer activity cell type specificity
1079
1080 # load libraries
1081 library(RColorBrewer)
1082 library(gplots)
1083 library(edgeR)
1084
1085 # load data
1086 enh <-
1087 read.table("STARR_starrs_peaks.enhancer_activity.eFDR05.ann.high_motif.scATAC_ACRs.b
1088 ed")
1089 acrs <- read.table("maize_scATAC_atlas_ACR_celltype_CPM.V5.txt")
1090 motifs <- read.table("TFBS_peaks.motifs.ENRICHED.enhancer_activity.bed")
1091
1092 # subset for representative leaf-derived clusters
1093 keep <- c("bulliform.2.26",
1094 "bundle_sheath.2.16",
```

```
1095 "ground_meristem.7.69",
1096 "guard_cell.7.74",
1097 "guard_mother_cell.7.71",
1098 "L1_SAM.4.46",
1099 "leaf_provascular.7.67",
1100 "mesophyll.2.14",
1101 "parenchyma.10.90",
1102 "protoderm.7.72",
1103 "stomatal_precursor.7.75",
1104 "subsidiary.7.68")
1105 all.acrs <- acrs
1106
1107 # rescale acrs
1108 acrs <- cpm(acrs, log=F)
1109 acrs <- acrs[,keep]
1110
1111 # subset enhancers by genomic feature
1112 enh <- subset(enh, enh$V8=="intergenic")
1113
1114 # get overlapping regions from the scATAC matrix
1115 enh$ids <- paste(enh$V11,enh$V12,enh$V13,sep="_")
1116 enh <- enh[order(enh$V11, decreasing=T),]
1117 enh <- enh[!duplicated(enh$ids),]
1118 shared <- intersect(enh$ids, rownames(acrs))
1119 rownames(enh) <- enh$ids
1120 enh <- enh[shared,]
1121 enh$starrIDs <- paste(enh$V1, enh$V2, enh$V3,sep="_")
1122
1123 # filter motifs
1124 motifs$starrIDs <- paste(motifs$V5, motifs$V6, motifs$V7, sep="_")
1125 motifs <- motifs[motifs$starrIDs %in% unique(enh$starrIDs),]
1126
1127 # normalize acrs
1128 acrs <- t(apply(acrs, 1, function(x){x/max(x)}))
1129
1130 # iterate over each cell type
1131 cts <- colnames(acrs)
1132 outs <- lapply(cts, function(x){
1133 access <- acrs[rownames(enh),x]
```



```
1134 names(access) <- enh$starrIDs
1135 motif.scores <- access[motifs$starrIDs] * as.numeric(as.character(motifs$V15))
1136 mtf <- data.frame(motif=motifs$V4, score=motif.scores)
1137 aves <- aggregate(score~motif, data=mtf, FUN=mean)
1138 score <- aves$score
1139 names(score) <- aves$motif
1140 return(score)
1141 })
1142 outs <- do.call(cbind, outs)
1143 colnames(outs) <- cts
1144 vars <- apply(outs, 1, var)
1145 outs <- outs[vars > 0,]
1146 z <- as.matrix(t(scale(t(outs))))
1147
1148 ## cluster columns
1149 co <- hclust(dist(t(outs)))$order
1150
1151 # reorder rows
1152 z <- z[,co]
1153 row.o <- apply(z, 1, which.max)
1154 z <- z[order(row.o, decreasing=F),]
1155
1156 # cap
1157 z[z < -3] <- -3
1158 z[z > 3] <- 3
1159
1160 # get family
1161 tfs <- data.frame(do.call(rbind, strsplit(rownames(z), "\\.")))
1162 cols2 <- colorRampPalette(brewer.pal(12, "Paired"))(length(unique(tfs$X1)))
1163 tfs$cols2 <- cols2[factor(tfs$X1)]
1164
1165 # visualize
1166 pdf("celltype_starr_motif_activity.pdf", width=10, height=10)
1167 heatmap.2(z, scale="none", trace='none',
1168 RowSideColors=tfs$cols,
1169 col=colorRampPalette(rev(brewer.pal(9, "RdBu")))(100),
1170 useRaster=T, Colv=F, Rowv=F, dendrogram="none", margins=c(9,9))
1171 dev.off()
1172
```

1173

1174 **Acknowledgments**

1175 This study was funded by support from the National Science Foundation (DBI-1906869) and the  
1176 National Institute of Health (1K99GM144742) to A.P.M. The ATAC-STARR-seq data analyzed in this  
1177 study was originally generated by Ricci, Lu, Ji and colleagues (Ricci *et al.*, 2019).

1178

1179

1180 **Competing interests**

1181 A.P.M. declares no competing interests.

1182

1183

1184 **Supplementary information**

1185 1. Data and code availability: All data and code have been deposited to GitHub:

1186 [https://github.com/Bio-protocol/Maize\\_ATAC\\_STARR\\_seq](https://github.com/Bio-protocol/Maize_ATAC_STARR_seq)

1187

1188

1189 **References**

1190 Arnold, C. D., Gerlach, D., Stelzer, C., Boryn, L. M., Rath, M. and Stark, A. (2013). Genome-Wide  
1191 Quantitative Enhancer Activity Maps Identified by STARR-seq. *Science* 339(6123): 1074-  
1192 1077. <Go to ISI>://WOS:000315452000041

1193 Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. and Greenleaf, W. J. (2013). Transposition of  
1194 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding  
1195 proteins and nucleosome position. *Nat Methods* 10(12): 1213-+. <Go to  
1196 ISI>://WOS:000327698100025

1197 Chen, S., Zhou, Y., Chen, Y. and Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ preprocessor.  
1198 *Bioinformatics* 34(17): i884-i890. <https://www.ncbi.nlm.nih.gov/pubmed/30423086>

1199 Grant, C. E., Bailey, T. L. and Noble, W. S. (2011). FIMO: scanning for occurrences of a given motif.  
1200 *Bioinformatics* 27(7): 1017-1018. <https://www.ncbi.nlm.nih.gov/pubmed/21330290>

1201 Hnisz, D., Abraham, B. J., Lee, T. I., Lau, A., Saint-Andre, V., Sigova, A. A., Hoke, H. A. and Young, R.  
1202 A. (2013). Super-enhancers in the control of cell identity and disease. *Cell* 155(4): 934-947.  
1203 <https://www.ncbi.nlm.nih.gov/pubmed/24119843>

1204 Hufford, M. B., Seetharam, A. S., Woodhouse, M. R., Chougule, K. M., Ou, S., Liu, J., Ricci, W. A.,  
1205 Guo, T., Olson, A., Qiu, Y., Della Coletta, R., Tittes, S., Hudson, A. I., Marand, A. P., Wei, S.,  
1206 Lu, Z., Wang, B., Tello-Ruiz, M. K., Piri, R. D., Wang, N., Kim, D. W., Zeng, Y., O'Connor, C.  
1207 H., Li, X., Gilbert, A. M., Baggs, E., Krasileva, K. V., Portwood, J. L., 2nd, Cannon, E. K. S.,  
1208 Andorf, C. M., Manchanda, N., Snodgrass, S. J., Hufnagel, D. E., Jiang, Q., Pedersen, S.,  
1209 Syring, M. L., Kudrna, D. A., Llaca, V., Fengler, K., Schmitz, R. J., Ross-Ibarra, J., Yu, J., Gent,  
1210 J. I., Hirsch, C. N., Ware, D. and Dawe, R. K. (2021). De novo assembly, annotation, and

- 1211 comparative analysis of 26 diverse maize genomes. *Science* 373(6555): 655-662.  
1212 <https://www.ncbi.nlm.nih.gov/pubmed/34353948>
- 1213 Jores, T., Tonnie, J., Dorrity, M. W., Cuperus, J. T., Fields, S. and Queitsch, C. (2020). Identification of  
1214 Plant Enhancers and Their Constituent Elements by STARR-seq in Tobacco Leaves[OPEN].  
1215 *The Plant Cell* 32(7): 2120-2131. <https://doi.org/10.1105/tpc.20.00155>
- 1216 Kent, W. J., Zweig, A. S., Barber, G., Hinrichs, A. S. and Karolchik, D. (2010). BigWig and BigBed:  
1217 enabling browsing of large distributed datasets. *Bioinformatics* 26(17): 2204-2207. <Go to  
1218 ISI>://WOS:000281738900023
- 1219 Leinonen, R., Sugawara, H., Shumway, M. and International Nucleotide Sequence Database, C.  
1220 (2011). The sequence read archive. *Nucleic Acids Res* 39(Database issue): D19-21.  
1221 <https://www.ncbi.nlm.nih.gov/pubmed/21062823>
- 1222 Li, H. (2011). Tabix: fast retrieval of sequence features from generic TAB-delimited files. *Bioinformatics*  
1223 27(5): 718-719. <https://www.ncbi.nlm.nih.gov/pubmed/21208982>
- 1224 Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform.  
1225 *Bioinformatics* 25(14): 1754-1760. <https://www.ncbi.nlm.nih.gov/pubmed/19451168>
- 1226 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin,  
1227 R. and Genome Project Data Processing, S. (2009). The Sequence Alignment/Map format and  
1228 SAMtools. *Bioinformatics* 25(16): 2078-2079. <https://www.ncbi.nlm.nih.gov/pubmed/19505943>
- 1229 Liu, T. (2014). Use model-based Analysis of ChIP-Seq (MACS) to analyze short reads generated by  
1230 sequencing protein-DNA interactions in embryonic stem cells. *Methods Mol Biol* 1150: 81-95.  
1231 <https://www.ncbi.nlm.nih.gov/pubmed/24743991>
- 1232 Marand, A. P., Chen, Z., Gallavotti, A. and Schmitz, R. J. (2021). A cis-regulatory atlas in maize at  
1233 single-cell resolution. *Cell* 184(11): 3041-3055 e3021.  
1234 <https://www.ncbi.nlm.nih.gov/pubmed/33964211>
- 1235 Marand, A. P., Zhang, T., Zhu, B. and Jiang, J. (2017). Towards genome-wide prediction and  
1236 characterization of enhancers in plants. *Biochim Biophys Acta Gene Regul Mech* 1860(1):  
1237 131-139. <https://www.ncbi.nlm.nih.gov/pubmed/27321818>
- 1238 Melnikov, A., Murugan, A., Zhang, X., Tesileanu, T., Wang, L., Rogov, P., Feizi, S., Gnirke, A., Callan,  
1239 C. G., Kinney, J. B., Kellis, M., Lander, E. S. and Mikkelsen, T. S. (2012). Systematic  
1240 dissection and optimization of inducible enhancers in human cells using a massively parallel  
1241 reporter assay. *Nature Biotechnology* 30(3): 271-277. <https://doi.org/10.1038/nbt.2137>
- 1242 Minnoye, L., Marinov, G. K., Krausgruber, T., Pan, L., Marand, A. P., Secchia, S., Greenleaf, W. J.,  
1243 Furlong, E. E. M., Zhao, K., Schmitz, R. J., Bock, C. and Aerts, S. (2021). Chromatin  
1244 accessibility profiling methods. *Nature Reviews Methods Primers* 1(1): 10.  
1245 <https://doi.org/10.1038/s43586-020-00008-9>
- 1246 Quinlan, A. R. and Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic  
1247 features. *Bioinformatics* 26(6): 841-842. <https://www.ncbi.nlm.nih.gov/pubmed/20110278>

- 1248 Ramirez, F., Dundar, F., Diehl, S., Gruning, B. A. and Manke, T. (2014). deepTools: a flexible platform  
1249 for exploring deep-sequencing data. *Nucleic Acids Res* 42(Web Server issue): W187-191.  
1250 <https://www.ncbi.nlm.nih.gov/pubmed/24799436>
- 1251 Ricci, W. A., Lu, Z., Ji, L., Marand, A. P., Ethridge, C. L., Murphy, N. G., Noshay, J. M., Galli, M., Mejia-  
1252 Guerra, M. K., Colome-Tatche, M., Johannes, F., Rowley, M. J., Corces, V. G., Zhai, J.,  
1253 Scanlon, M. J., Buckler, E. S., Gallavotti, A., Springer, N. M., Schmitz, R. J. and Zhang, X.  
1254 (2019). Widespread long-range cis-regulatory elements in the maize genome. *Nat Plants*  
1255 5(12): 1237-1249. <https://www.ncbi.nlm.nih.gov/pubmed/31740773>
- 1256 Schmitz, R. J., Grotewold, E. and Stam, M. (2022). Cis-regulatory sequences in plants: Their  
1257 importance, discovery, and future challenges. *Plant Cell* 34(2): 718-741. <Go to  
1258 ISI>://WOS:000761460300008
- 1259 Sun, J., He, N., Niu, L., Huang, Y., Shen, W., Zhang, Y., Li, L. and Hou, C. (2019). Global Quantitative  
1260 Mapping of Enhancers in Rice by STARR-seq. *Genomics Proteomics Bioinformatics* 17(2):  
1261 140-153. <https://www.ncbi.nlm.nih.gov/pubmed/31201999>
- 1262 Thorvaldsdottir, H., Robinson, J. T. and Mesirov, J. P. (2013). Integrative Genomics Viewer (IGV): high-  
1263 performance genomics data visualization and exploration. *Brief Bioinform* 14(2): 178-192.  
1264 <https://www.ncbi.nlm.nih.gov/pubmed/22517427>
- 1265 Zhao, H., Sun, Z. F., Wang, J., Huang, H. J., Kocher, J. P. and Wang, L. G. (2014). CrossMap: a  
1266 versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* 30(7):  
1267 1006-1007. <Go to ISI>://WOS:000334078300017  
1268  
1269