

# End Motif Tutorial

Jeffrey Szymanski

August 23, 2024

Minimal working example.

Dependencies

- R and python packages as specified per script.
- Small bam files human1-5 in data/bams
- Appropriate reference fasta, GCA\_000001405.15\_GRCh38\_no\_alt\_analysis\_set.fna (available from [ftp://ftp.ncbi.nlm.nih.gov/genomes/UCSC/Genomes/Human/1405.15/GRCh38/Genomes/Homo\\_sapiens/GCA\\_000001405.15\\_GRCh38\\_no\\_alt\\_analysis\\_set.fna.gz](ftp://ftp.ncbi.nlm.nih.gov/genomes/UCSC/Genomes/Human/1405.15/GRCh38/Genomes/Homo_sapiens/GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.gz))

## 1. Clear motifs data directory if present:

```
1 if [ -d data/motifs ]; then rm -rf data/motifs; fi
2 mkdir data/motifs
```

## 2. Generate motifs from 5' ends of bams:

```
1 scripts/sample_motifs.py -h
2
3 REFERENCE_GENOME="/home/jeszyman/pnst/inputs/GCA_000001405.15_GRCh38_no_alt_analysis_set.fna"
4
5 scripts/sample_motifs.py --bam_file data/bams/human1.bam \
6     --output_file data/motifs/human1_motifs.tsv \
7     --motif_length 4 \
8     --num_reads 10000 \
9     --threads 4 \
10    --reference_genome $REFERENCE_GENOME
11
12 for bam in data/bams/*.bam; do
13     base_name=$(basename "$bam" .bam)
14
15     output_file="data/motifs/${base_name}_motifs.tsv"
16
17     scripts/sample_motifs.py --bam_file "$bam" --output_file "$output_file" \
18         --motif_length 4 --num_reads 10000 \
19         --threads 4 --reference_genome "$REFERENCE_GENOME"
20 done
```

### 3. Consolidate to a single matrix:

(In R)

```
1 library(tidyverse)
2
3 # Function to read a file and format it for merging with counts
4 read_motif_file_counts <- function(file) {
5   df <- read_tsv(file, col_names = c("motif", "count"))
6   file_name <- tools::file_path_sans_ext(basename(file))
7   df <- df %>% rename(!!file_name := count)
8   return(df)
9 }
10
11 # Function to read a file and format it for merging with fractions
12 read_motif_file_fractions <- function(file) {
13   df <- read_tsv(file, col_names = c("motif", "count"))
14   total_count <- sum(df$count)
15   df <- df %>% mutate(fraction = count / total_count)
16   file_name <- tools::file_path_sans_ext(basename(file))
17   df <- df %>% select(motif, fraction) %>% rename(!!file_name := fraction)
18   return(df)
19 }
20
21 # List of files
22 files <- list.files(path = "data/motifs", pattern = "*.tsv", full.names = TRUE)
23
24 # Read and merge all files for counts
25 motif_data_counts <- files %>%
26   map(read_motif_file_counts) %>%
27   reduce(full_join, by = "motif")
28
29 # Read and merge all files for fractions
30 motif_data_fractions <- files %>%
31   map(read_motif_file_fractions) %>%
32   reduce(full_join, by = "motif")
33
34 # Save the resulting matrices
35 write_tsv(motif_data_counts, "data/combined_motif_counts_matrix.tsv")
36 write_tsv(motif_data_fractions, "data/combined_motif_fractions_matrix.tsv")
```

### 4. Motif diversity score

(In R)

---

```

1 annotation = data.frame(library =
  ↪ c("human1_motifs", "human2_motifs", "human3_motifs", "human4_motifs", "human5_motifs"),
2       cohort = c("healthy", "cancer", "healthy", "cancer", "healthy"))
3 annotation
4
5 motifs = read_tsv("data/combined_motif_fractions_matrix.tsv")
6
7 motifs
8
9 motifs_long <-
10   pivot_longer(motifs, cols = !motif, names_to = "library", values_to = "fraction") %>%
11   left_join(annotation, by = "library") %>%
12   select(motif, library, fraction, cohort)
13
14 motifs_long
15
16 mds = motifs_long %>%
17   mutate(fraction = ifelse(fraction == 0, 1e-10, fraction)) %>% # Add a small constant
  ↪ to avoid log(0)
18   mutate(mds = -fraction * log(fraction) / log(256)) %>%
19   group_by(library) %>%
20   summarize(mds = sum(mds))
21
22 mds

```

---

## 5. F-profiles

---

```

1 scripts/fprofiles.py -h
2
3 scripts/fprofiles.py --data_file ./data/combined_motif_fractions_matrix.tsv
  ↪ --output_fprof_per_lib ./data/fprofiles.tsv --output_motif_per_fprof
  ↪ ./data/motif_per_fprofile.tsv

```

---