

Cheat sheet for the users parameters, and output file structures of CSREP

This cheat sheet aims at helping you easily follow along the tutorial, and later easily modify the config/config.yml for your own usage of CSREP or base_count method.

Files/folders that contain CSREP's output (or supporting files produced by the pipeline)

Files/folders that users provide

- **Bold, black and italic words** represent variable names as appeared in config/config.yml
- **Bold and black words** show the example file path as presented in the testdata folder provided along with the tutorial

raw_user_input_dir

Where users store all input segmentation data files for all samples, regardless of their group memberships. **testdata/raw_data**

<sample_id> Ex: E003

Each of these folder corresponds to one sample. **testdata/raw_data/E003/**

<sampleID><input_filename_suffix>

Input segmentation data for sample **<sampleID>**. **<input_filename_suffix>** is shared across all samples, and is a user input param. **testdata/raw_data/E003/E003_chr22_core_K27ac_segments.bed.gz** → **_chr22_core_K27ac_segments.bed.gz** correspond to **input_filename_suffix** in **config/config.yml**

state_annot_fn

This file specifies different characteristics of the states . Two mandatory columns: **state** and **mnemonic**. **testdata/state_annot.txt**

chrom_length_fn

bed file showing the length of each chromosome, used when CSREP try to generate random sample regions across the genome for training. **testdata/roadmap_18state_chromlength.bed**

sample_genome_fn

where the data of regions that we sample for training data are stored (we recommend specifying this file path inside **all_ct_out_dir**). **testdata/sample_genome.bed.gz**

training_data_folder

the folder where we will store chromatin state segmentation data used for training, which corresponds to 10% of the genome, for each sample.

