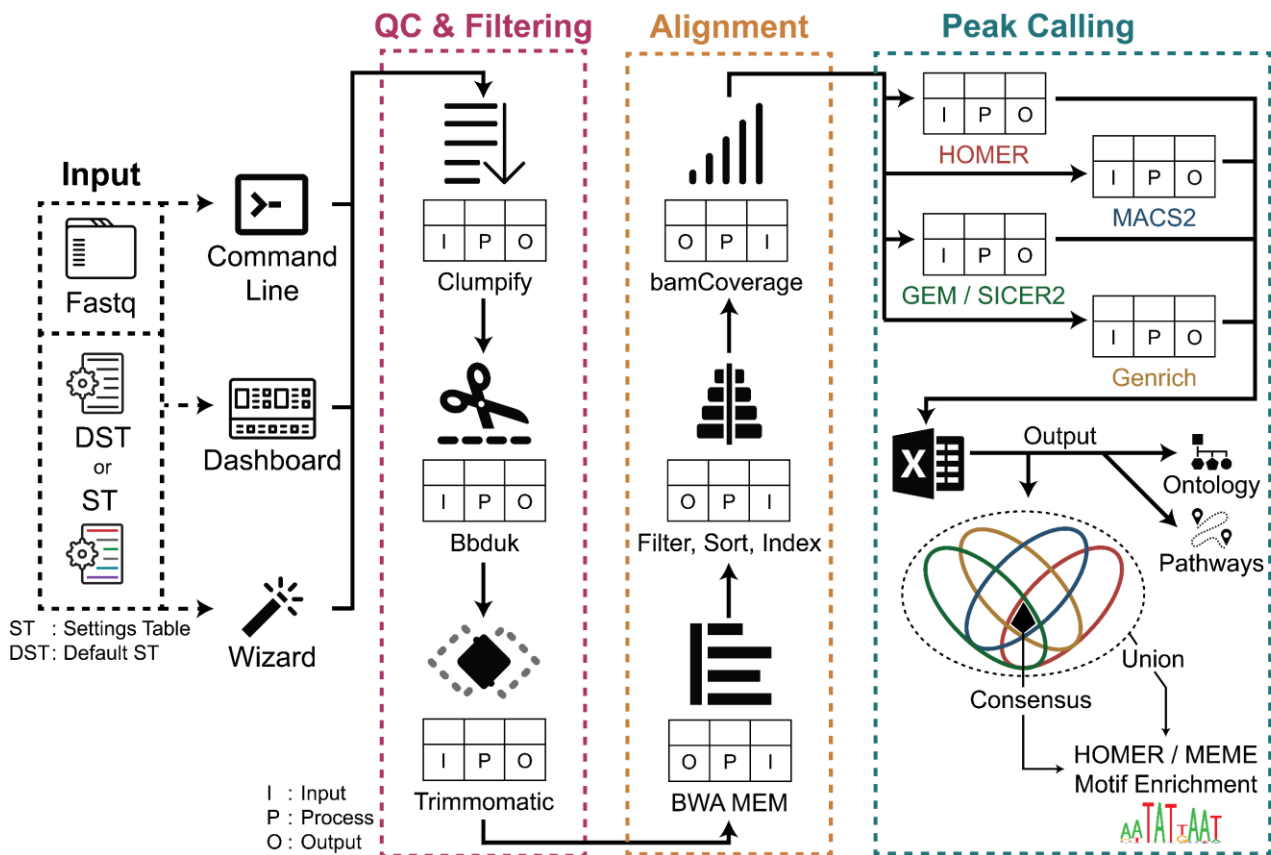




## ChIP-AP Graphical Overview



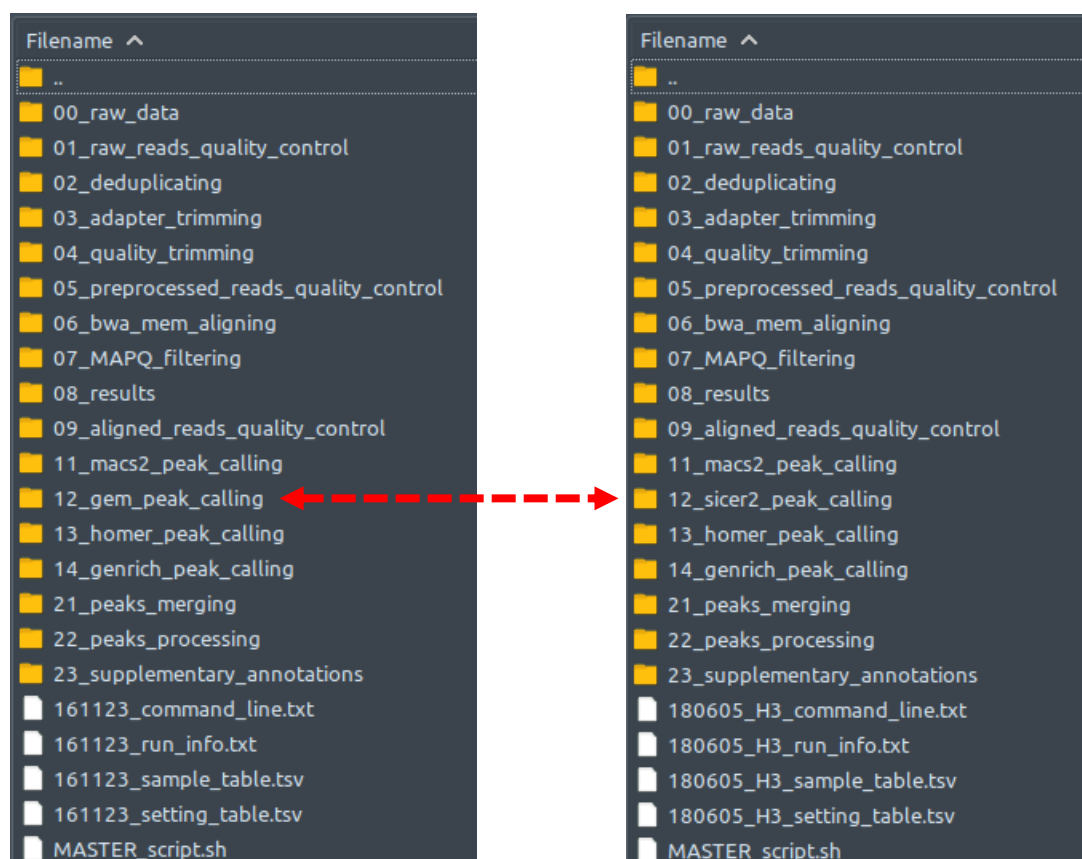
## Detailed Explanation of Steps and Methodology Used

- 1. Acquisition of raw sequencing files.** ChIP-AP can directly process the output files of sequencing instruments. Files may be in FASTQ, or compressed FASTQ (.gz) format. ChIP-AP can also process aligned reads in BAM format. Reads may be single or paired ends. Background control is compulsory, no unmatched samples allowed here.
- 2. Sample recognition and registration.** Performed by the main script. Each input sample is registered into the system and given a new name according to their sample category (ChIP or background control), replicate number, and whether it's the first or second read file (in case of paired end sequencing data). Afterwards, their formats and compression status is recognized and processed into gun-zipped FASTQ as necessary.



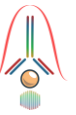
3. **Generation of multiple modular scripts.** Each process in ChIP-AP is executed from individual scripts generated. This was an intentional design decision as it allows for easy access for modifications of any step within the pipeline without hunting through the master ChIP-AP scripts. You can thank us later if you have to modify and tailor something later and don't have to drudge through the trenches of someone else's code. Chocolate treats always welcome!

**How does it look like?** After this step is done, which is practically the end of your ChIP-AP processes if you don't use the `--run` flag to run the pipeline immediately, you can see within your designated output directory a single folder named based on your `--setname` input, with contents just like below regardless of single-end or paired end-mode:



The left figure shows a dataset with narrow peak type. The right figure shows a dataset with broad peak type. Note that the peak calling module number 12 is interchangeable between GEM for narrow peak type and SICER2 for broad peak type.

Each of these folders are basically empty, and contains a script which is named based on the folder name (e.g., script **02\_deduplicating.sh** inside folder **02\_deduplicating**). Each of these scripts will be executed in numerical sequence when you run the pipeline. Aside from these scripts, there will be several miscellaneous text files which contain essential information of your pipeline run (See **Miscellaneous Pipeline Output** section below for details). Lastly, there is your big red button: the **MASTER\_script.sh** that you can simply call to sequentially run all the scripts within the aforementioned folders.

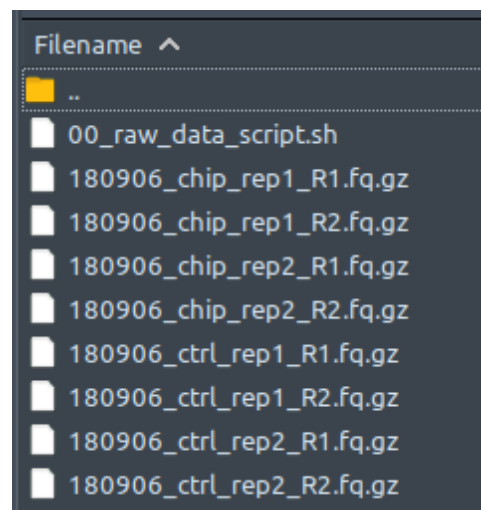
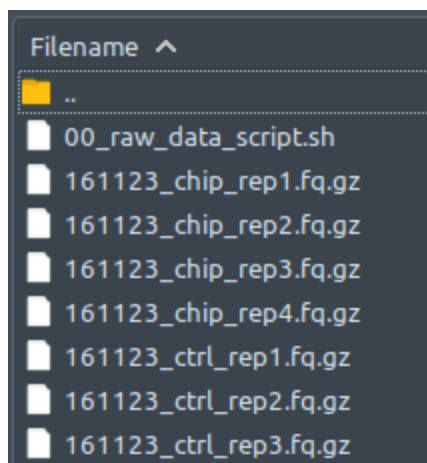


4. **Copying, compressing, and renaming of the raw sequencing reads.** In the very beginning, ChIP-AP makes (in the user-designated output folder) a copy of each unaligned sequence reads file (e.g., fastq), compresses them into a gunzipped file (if not already), and renames them with the prepared new name from step “2. **Sample recognition and registration**”. If the given inputs are aligned reads (bam files), the pipeline starts at step “12. **Sorting and indexing of aligned reads files**” (see below) and the copying and renaming are taken over by “08\_results\_script.sh” where the original bam files are directly sorted and the pipeline proceeds normally from there.

Modular script used: **00\_raw\_data\_script.sh**

- **Operation:** cp, gzip, mv (Bash)
- **Input** : [origin directory] / [original ChIP/control filename]
- **Process** : Copy, compress, and rename raw reads files
- **Output** : [output directory] / 00\_raw\_data / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].fq.gz

**How does it look like?** After **00\_raw\_data\_script.sh** had been executed, every single reads file in your dataset will be copied into this folder: **00\_raw\_data**, compressed into **fq.gz**, and renamed into something like the preview below, regardless of your initial filenames, fastq formatting or extensions.



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). All these **fq.gz** files will be immediately deleted at the end of 02\_deduplicating\_script.sh execution if --deltemp flag is used on ChIP-AP call. We won't explain how to open and read these **fq.gz** files, since if you need us to tell you that, you most probably will not be able to make anything out of the contents in there anyway.

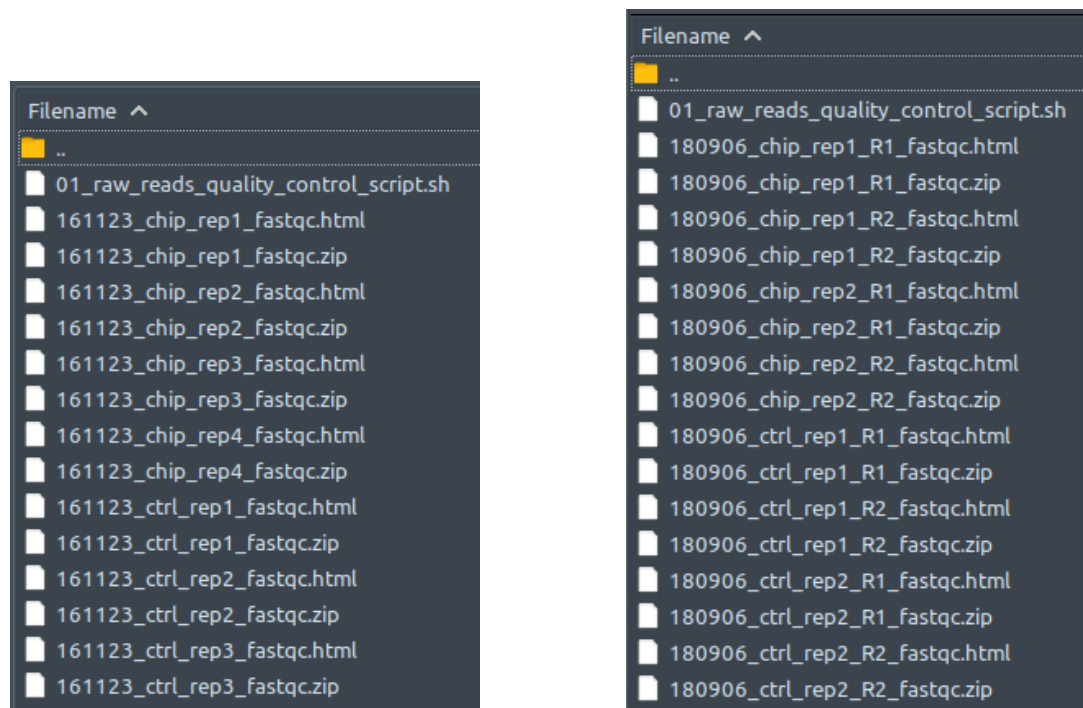


5. **Raw sequencing reads quality assessment.** Performed by FastQC. Reads quality assessment is performed to check for duplicates, adapter sequences, base call scores, etc. Assessment results are saved as reports for user viewing. If the final results are not as expected, it's worthwhile to go through the multiple QC steps and track the quality of the data as it's processed. If the default QC steps aren't cleaning up the data adequately, you may need to modify some parameters to be more/less stringent with cleanup. From our testing, our default values seem to do a fairly adequate job though for most datasets.

Modular script used: **01\_raw\_reads\_quality\_control\_script.sh**

- **Calls** : fastqc
- **Input** : 00\_raw\_data / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].fq.gz
- **Process** : Generate raw reads quality assessment reports
- **Output** : 01\_raw\_reads\_quality\_control / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2]\_fastqc.html

**How does it look like?** After **01\_raw\_reads\_quality\_control\_script.sh** has been executed, the folder: **01\_raw\_reads\_quality\_control** will contain all these quality assessment reports for every raw reads file in folder **00\_raw\_data**, just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). The **.zip** files contain the individual components to be compiled for the report so you can ignore those. To read the reports, open the **.html** files. This file is a multitabular file in which you can evaluate the quality of your experiment, sequencing, etc. Comprehensive as it is, explaining the contents in detail would take a whole new guide by itself. Therefore, in case the reports do not spell everything out enough for you, check this out: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/>.

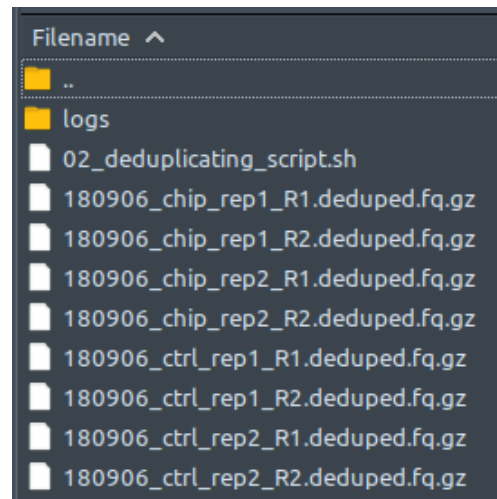
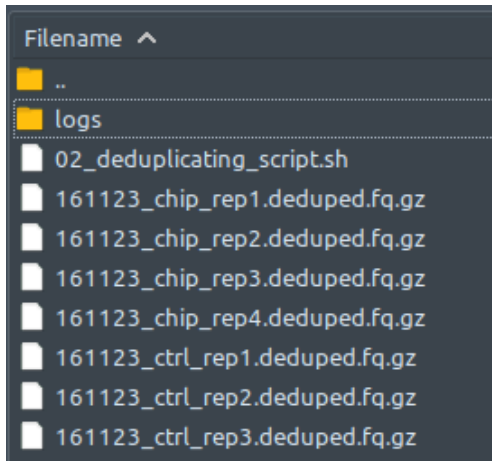


6. **Deduplication of reads.** Performed by clumpify from BBMap package. Necessary command line argument is given to clumpify in order to remove optical duplicates and tile-edge duplicates from the reads file in addition to PCR duplicates. Optimization of file compression is also performed by clumpify during deduplication process, in order to minimize storage space and speed up reads file processing.

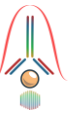
Modular script used: **02\_deduplicating\_script.sh**

- **Calls** : clumpify.sh
- **Input** : 00\_raw\_data / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].fq.gz
- **Process** : Remove PCR duplicates, optical duplicates, and tile-edge duplicates
- **Output** : 02\_deduplicating / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].deduped.fq.gz

**How does it look like?** After **02\_deduplicating\_script.sh** had been executed, the folder: **02\_deduplicating** will contain all these deduplicated reads files (marked by the extension: **.deduped.fq.gz**), just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). There should be one deduplicated file for each processed raw reads file from folder **00\_raw\_data**. Paired-end files are processed in pairs by **clumpify**. All these **deduped.fq.gz** files will be deleted at the end of **02\_deduplicating\_script.sh** execution if **--deltemp** flag is used on ChIP-AP call.

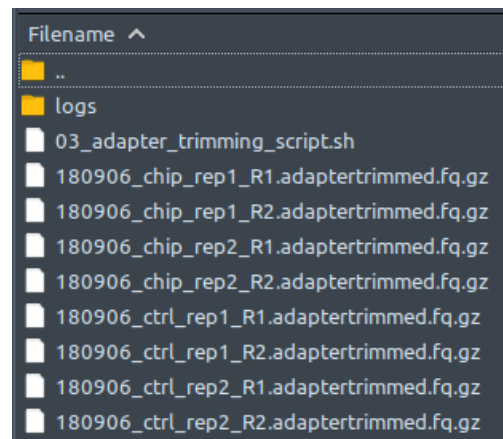
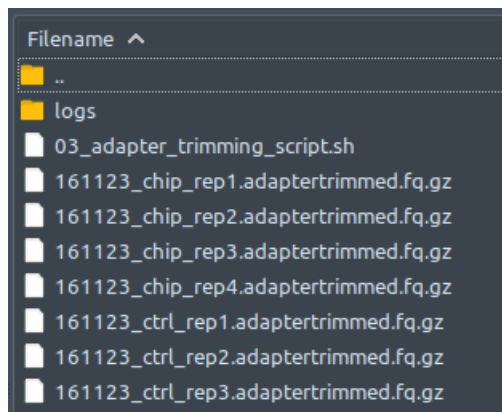


7. **Adapter trimming of reads.** Performed by BBDuk from BBMap package. BBDuk scans every read for adapter sequence, based on the reference list adapters given in the command line argument. The standard BBDuk adapter sequence reference list 'adapter.fa' is used as a default in the pipeline. Any sequencing adapter present in the reads is removed. Custom adapter sequence can be used whenever necessary or by modifying the adapter.fa file with your new sequences.

Modular script used: **03\_adapter\_trimming\_script.sh**

- **Calls** : bbdduk.sh
- **Input** : 02\_deduplicating / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].deduped.fq.gz  
[path to genome folder] / bbmap / adapters.fa (file provided by ChIP-AP)
- **Process** : Trim away adapter sequences based on given sequences in file 'adapters.fa'
- **Output** : 03\_adapter\_trimming / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].adaptertrimmed.fq.gz

**How does it look like?** After **03\_adapter\_trimming\_script.sh** had been executed, the folder: **03\_adapter\_trimming** will contain all these adapter-trimmed reads files (marked by the extension: **.adaptertrimmed.fq.gz**), just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). There should be one adapter-trimmed file for each processed deduplicated reads file from folder **02\_deduplicating**. Paired-end files are processed in pairs by **bbduk**. All these **adaptertrimmed.fq.gz** files will be deleted at the end of **03\_adapter\_trimming\_script.sh** execution if **--deltemp** flag is used on ChIP-AP call.



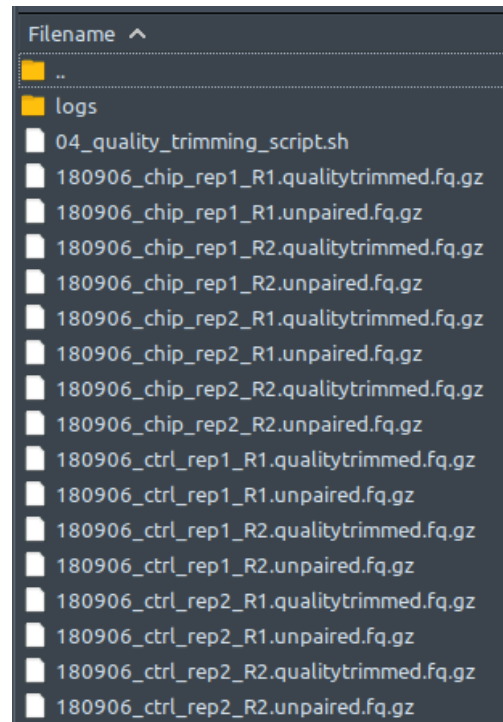
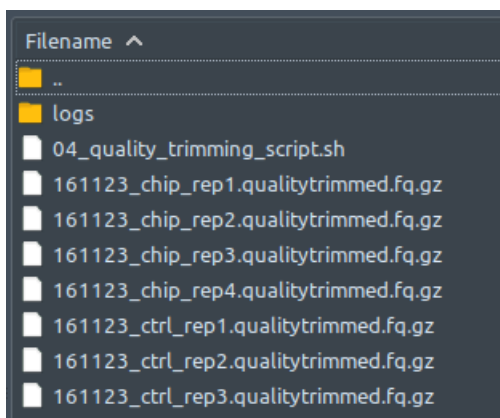


8. **Quality trimming of reads.** Performed by **trimmomatic**. Trimmomatic scans every read trims low quality base calls from reads. Additionally, it scans with a moving window along the read and cuts the remainder of the read when the average quality of base calls within the scanning window drops below the set threshold. Finally, it discards the entirety of a read if it gets too short post-trimming for alignment to reference genome, minimizing the chance of reads being multi-mapped to multiple genomic locations.

Modular script used: **04\_quality\_trimming\_script.sh**

- **Calls** : trimmomatic
- **Input** : 03\_adapter\_trimming / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].adaptertrimmed.fq.gz
- **Process** : Remove reads with low PHRED (base calling) score
- **Output** : 04\_quality\_trimming / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].qualitytrimmed.fq.gz

**How does it look like?** After **04\_quality\_trimming\_script.sh** had been executed, the folder: **04\_quality\_trimming** will contain all these quality-trimmed reads files (marked by the extension: **.qualitytrimmed.fq.gz**), just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). There should be one quality-trimmed file for each processed adapter-trimmed reads file from folder **03\_adapter\_trimming**. Paired-end files are processed in pairs by **trimmomatic**. Unpaired reads are separated (saved into **unpaired.fq.gz** files) from the paired reads (saved into **qualitytrimmed.fq.gz** files). Only the paired reads (extension: **.qualitytrimmed.fq.gz**) are processed further in the pipeline. All these **qualitytrimmed.fq.gz** and **unpaired.fq.gz** files will be immediately deleted at the end of **04\_quality\_trimming\_script.sh** execution if **--deltemp** flag is used on ChIP-AP call.

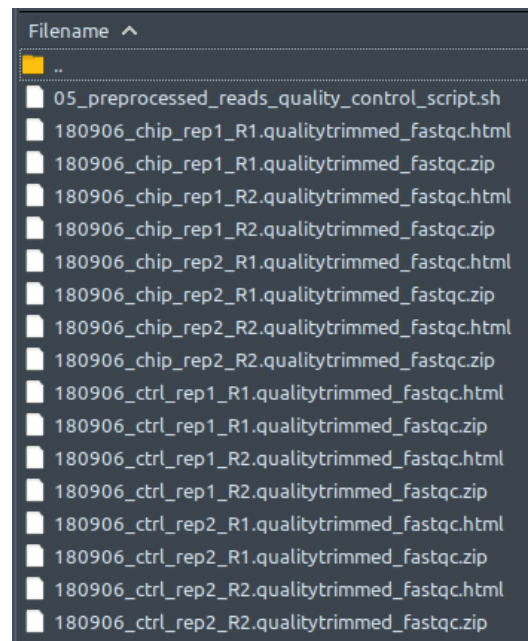
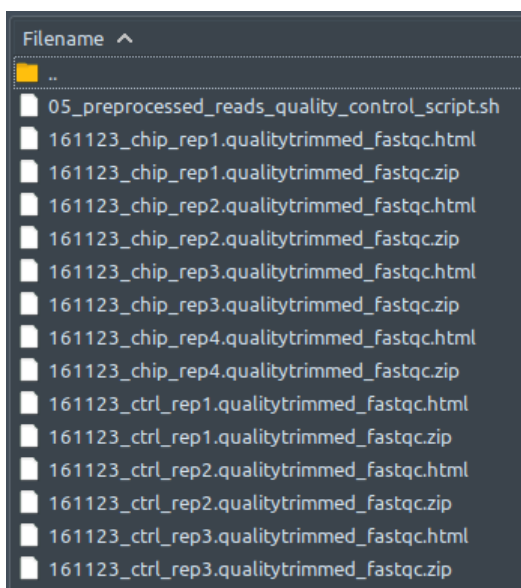


9. **Pre-processed reads quality assessment.** Performed by FastQC. Quality assessment is performed to check for the efficiency of cleanup. Results are saved as reports.

Modular script used: **05\_preprocessed\_reads\_quality\_control\_script.sh**

- **Calls** : fastqc
- **Input** : 04\_quality\_trimming / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].qualitytrimmed.fq.gz
- **Process** : Generate preprocessed reads quality assessment reports
- **Output** : 05\_preprocessed\_reads\_quality\_control / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2]\_fastqc.html

**How does it look like?** After **05\_preprocessed\_reads\_quality\_control\_script.sh** had been executed, the folder: **05\_preprocessed\_reads\_quality\_control** will contain all these quality assessment reports for every **qualitytrimmed.fq.gz** file in folder **04\_quality\_trimming**, just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples, in which every sample consists of two files: the first read (R1), and the second read (R2). The **.zip** files contains the individual components to be compiled for the report so you can ignore those. To read the reports, open the **.html** files. This file is a multitabular file in which you can evaluate how your preprocessings: deduplication, adapter trimming, and quality trimming, affected your sequencing reads. Comprehensive as it is, explaining the contents in detail would take a whole new guide by itself. Therefore, in case the reports do not spell everything out enough for you, check this out: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/>.



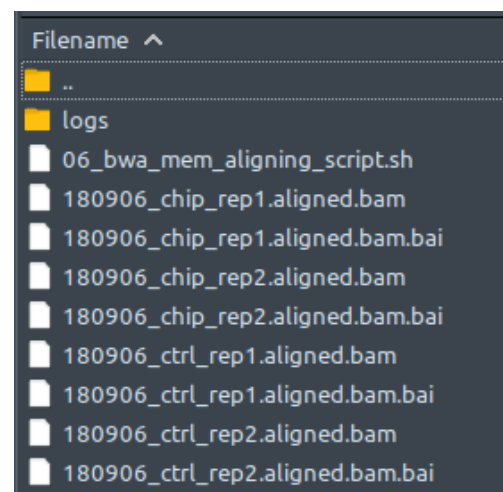
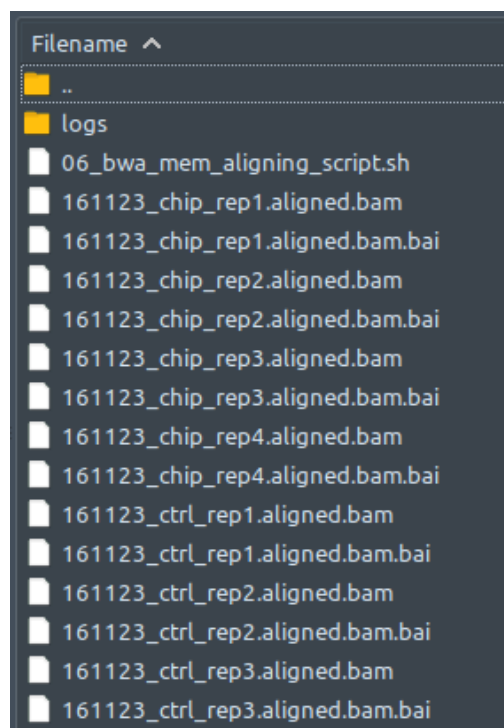


10. **Reads alignment to reference genome.** Performed by the mem algorithm in BWA aligner. Appropriate genome reference for the sample organism is given as a command line argument. The default genome reference is hg38. Precomputed genome references hg38, hg19, mm9, mm10, dm6, and sacCer3 are downloaded as part of the ChIP-AP installation process. Tutorials for custom/different genome references can be found on the ChIP-AP github (coming soon!). BWA is used in preference to other aligners such as Bowtie2 as in benchmarking papers (Thankaswamy-Kosalai et al., 2017), we were more satisfied with the results of BWA, hence its inclusion in this pipeline.

Modular script used: **06\_bwa\_mem\_aligning\_script.sh**

- **Calls** : bwa mem
- **Input** : 04\_quality\_trimming / [setname]\_[chip/ctrl]\_rep[#]\_R[1/2].qualitytrimmed.fq.gz  
[path to genome folder] / bwa / (reference genome provided by ChIP-AP)
- **Process** : Align preprocessed reads to the designated reference genome
- **Output** : 06\_bwa\_mem\_aligning / [setname]\_[chip/ctrl]\_rep[#].aligned.bam

**How does it look like?** After **06\_bwa\_mem\_aligning\_script.sh** had been executed, the folder: **06\_bwa\_mem\_aligning** will contain all these aligned reads files (marked by the extension: **.aligned.bam**), just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples. Note that right here the first reads (R1), and the second reads (R2) had both been aligned into the same reference genome, and thus no longer separated in two different files. Paired-end files are processed in pairs by **bwa mem**. There should be one aligned reads file here for each processed single-end, or for every two processed paired-end quality-trimmed reads file from folder **04\_quality\_trimming**. All these **aligned.bam** files will be deleted at the end of **06\_bwa\_mem\_aligning\_script.sh** execution if **--deltemp** flag is used on ChIP-AP call.

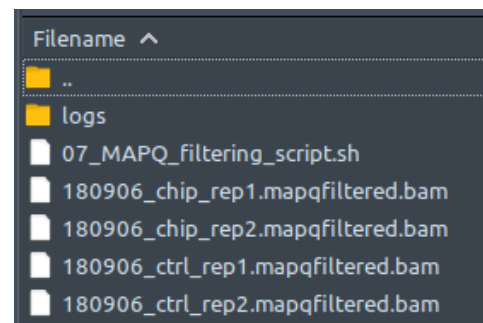
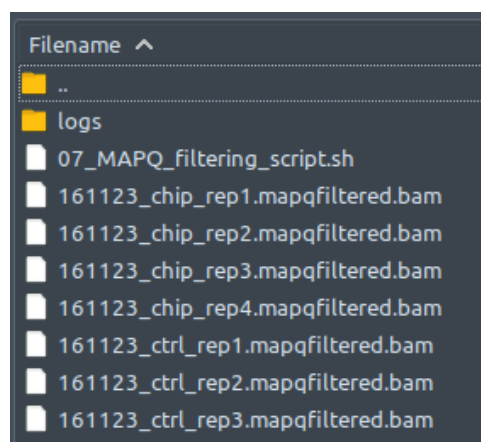


11. **Alignment score quality filtering.** Performed by samtools view. This filter (if set) will remove all reads with alignment score (MAPQ) below a user defined threshold. Reads with suboptimal fit into the genome and/or reads with multiple ambiguous mapped locations can easily be excluded from the reads file using this filter step also. To disable MAPQ filtering, simply remove all flags from the settings table for this step.

Modular script used: **07\_MAPQ\_filtering\_script.sh**

- **Calls** : samtools view
- **Input** : 06\_bwa\_mem\_aligning / [setname]\_[chip/ctrl]\_rep[#].aligned.bam
- **Process** : Remove reads with low MAPQ (alignment) score
- **Output** : 07\_MAPQ\_filtering / [setname]\_[chip/ctrl]\_rep[#].mapqfiltered.bam

**How does it look like?** After **07\_MAPQ\_filtering\_script.sh** had been executed, the folder: **07\_MAPQ\_filtering** will contain all these MAPQ-filtered reads files (marked by the extension: **.mapqfiltered.bam**), just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples. Again, note that right here the first reads (R1), and the second reads (R2) had both been aligned into the same reference genome by bwa mem in above, and thus no longer separated in two different files. There should be one MAPQ-filtered reads file here for each processed aligned reads file from folder **06\_bwa\_mem\_aligning**. All these **mapqfiltered.bam** files will be deleted at the end of **07\_MAPQ\_filtering\_script.sh** execution if --deltemp flag is used on CHIP-AP call.



12. **Sorting and indexing of aligned reads files.** Performed by samtools sort and samtools index, which do nothing to the aligned reads files other than sorting and indexing, priming the aligned reads files for further processing.

Modular script used: **08\_results\_script.sh**

- **Calls** : samtools sort
- **Input** : 07\_MAPQ\_filtering / [setname]\_chip/ctrl\_rep[#].mapqfiltered.bam
- **Process** : Sort all bam files based on coordinate
- **Output** : 08\_results / [setname]\_chip/ctrl\_rep[#].bam
  
- **Calls** : samtools merge
- **Input** : 08\_results / [setname]\_chip\_rep[#].bam  
08\_results / [setname]\_ctrl\_rep[#].bam
- **Process** : Merge all sorted ChIP bam files and all sorted control bam files.
- **Output** : 08\_results / [setname]\_chip\_merged.bam  
08\_results / [setname]\_ctrl\_merged.bam
- **Condition:** *--fcmerge flag is used OR unequal number of ChIP and control samples  
OR peak type is broad*
  
- **Calls** : samtools index
- **Input** : 08\_results / [setname]\_chip/ctrl\_rep[#].bam  
08\_results / [setname]\_chip/ctrl\_merged.bam
- **Process** : Make indices for all coordinate-sorted bam files
- **Output** : 08\_results / [setname]\_chip/ctrl\_rep[#].bam.bai  
08\_results / [setname]\_chip/ctrl\_merged.bam.bai
  
- **Calls** : samtools sort -n
- **Input** : 08\_results / [setname]\_chip/ctrl\_rep[#].bam
- **Process** : Sort all bam files based on read name
- **Output** : 08\_results / [setname]\_chip/ctrl\_rep[#]\_namesorted.bam

**How does it look like?** Detailed output preview of step 12, 13, and 14 are combined below step 14

13. **ChIP pulldown efficiency assessment.** Performed by plotFingerprint from the deeptools package, which generates fingerprint plots. These serve as a quality control figure that shows DNA pulldown efficiency of the ChIP experiment. Refer to the appropriate documentation for full details but in short – the input should be as close to the 1:1 diagonal as possible and the better enrichment seen in your sample, the more its curve will bend towards the bottom right. You want (ideally) a large gap between the chip and the control samples. PNG files are provided for easy viewing, SVG files provided if you want to make HQ versions later for publication.

Modular script used: **08\_results\_script.sh**

- **Calls** : plotFingerprint
- **Input** : 08\_results / [setname]\_chip/ctrl\_rep[#].bam
- **Process** : Generate fingerprint plots for all bam files
- **Output** : 08\_results / fingerprint\_plots/[setname].[png/svg]  
08\_results / fingerprint\_plots/[setname]\_merged.[png/svg]

**How does it look like?** Detailed output preview of step 12, 13, and 14 are combined below step 14



14. **Visualization track generation of aligned reads files.** Performed by bamCoverage from the deeptools package. Generates bigwig files for quick and simple visualization of reads distribution along the referenced genome using local tools such as IGV. The Coverage tracks can be uploaded to genome browsers such as UCSC, however a track hub needs to be generated – something ChIP-AP does not do at this stage.

Modular script used: **08\_results\_script.sh**

- **Calls** : bamCoverage
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bam  
08\_results / [setname]\_[chip/ctrl]\_merged.bam
- **Process** : Generate BigWig coverage file for each individual bam file
- **Output** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bw  
08\_results / [setname]\_[chip/ctrl]\_merged.bw

**How does it look like?** After **08\_results\_script.sh** had been executed, the folder: **08\_results** will contain all these sorted reads files (marked by the extension: **.bam**), a couple merged sorted reads files\* (marked by the extension: **\_merged.bam**), indices to all the sorted reads files (marked by the extension: **.bam.bai**), the same sorted reads files re-sorted by name (marked by the extension: **namesorted.bam**), bigwig files of all the sorted reads files (marked by the extension: **.bw**), fingerprint plot files of all the sorted reads files (in its own folder: **fingerprint\_plots**), and all the log files from all the program calls by **08\_results\_script.sh**.

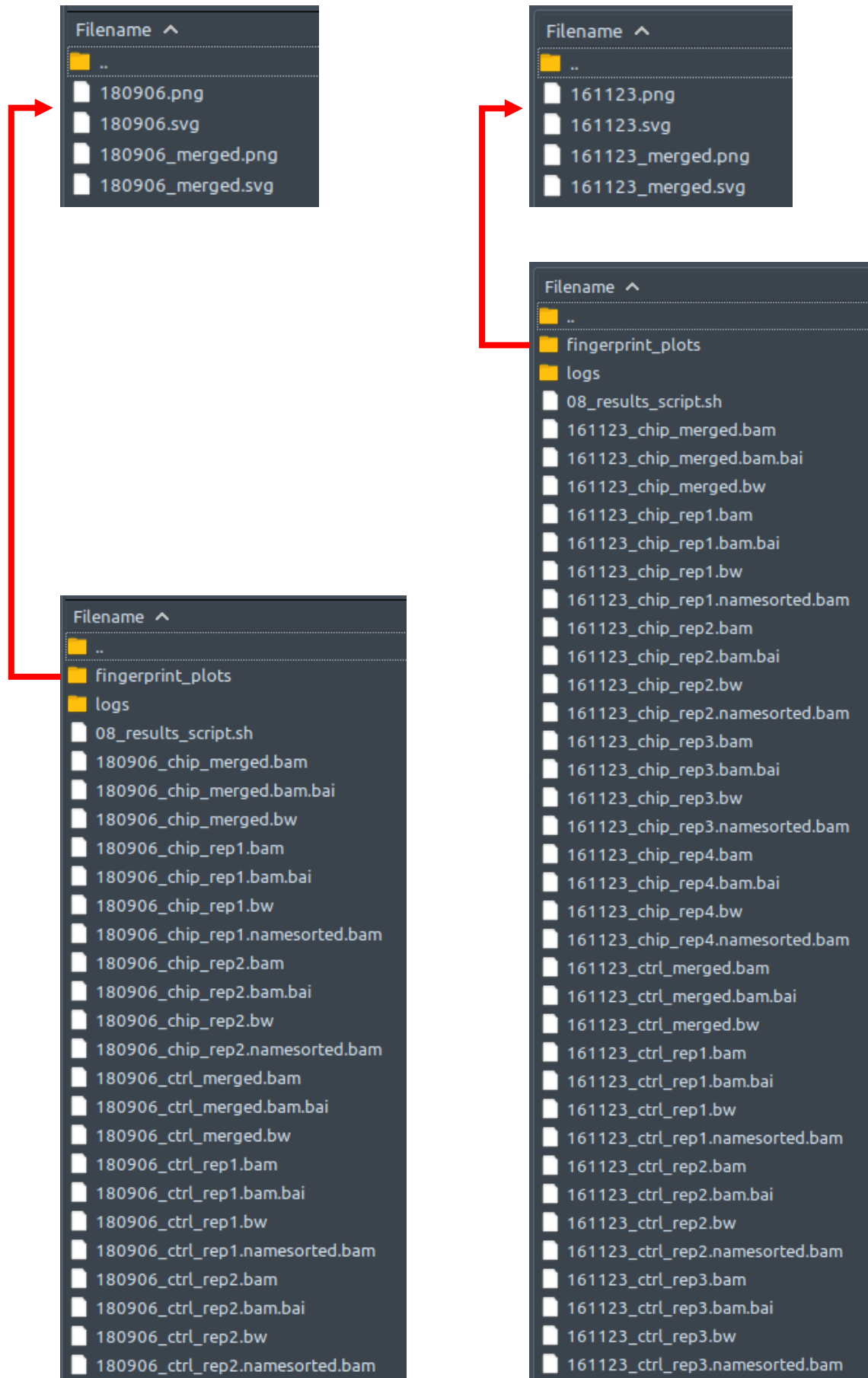
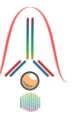
The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples. There should be one sorted reads file here for each processed MAPQ-filtered reads file from folder **07\_MAPQ\_filtering**; a merged ChIP and a merged control reads files\*; one index file for each sorted reads file in this folder; one name-sorted reads file for each sorted reads file in this folder\*\*; one bigwig file for each sorted reads file in this folder; two fingerprint plot files in the **fingerprint\_plots** folder (in extension: **.png** and **.svg**); and another two fingerprint plot files in the **fingerprint\_plots** folder (in extension: **.png** and **.svg**)\*\*\*.

To view and analyze the read distribution of your peaks, load the BigWig files (**.bw**) into the the program: **IGV** (more details down below). To view and evaluate your ChIP experiment DNA pulldown efficiency, open the **.png** or **.svg** using any supporting image viewer program (more details down below).

\* Only when needed by the pipeline

\*\* Except the merged ones

\*\*\* Only when merged sorted reads files are present



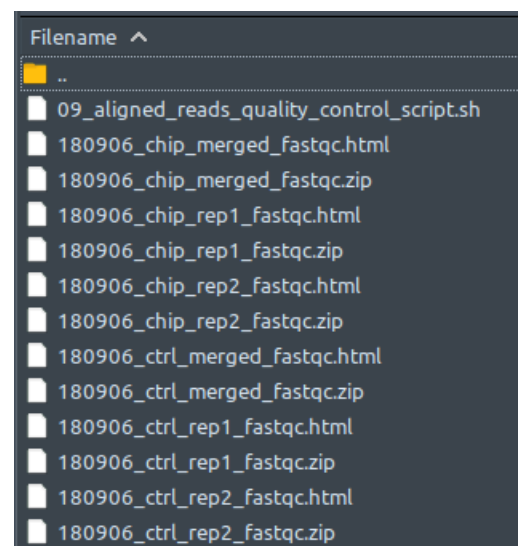
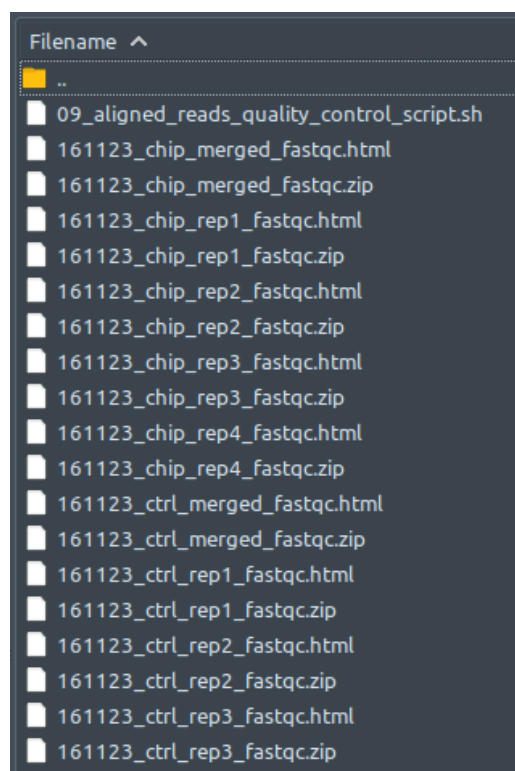


15. **Aligned reads quality assessment.** Processed by FastQC. Quality assessment is performed to check for the alignment efficiency, such as how many reads failed to be mapped. Assessment results are saved as reports.

Modular script used: **09\_aligned\_reads\_quality\_control\_script.sh**

- **Calls** : fastqc
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bam  
08\_results / [setname]\_[chip/ctrl]\_merged.bam
- **Process** : Generate raw reads quality assessment reports
- **Output** : 08\_results / [setname]\_[chip/ctrl]\_rep[#]\_fastqc.html  
08\_results / [setname]\_[chip/ctrl]\_merged\_fastqc.html

**How does it look like?** After **09\_aligned\_reads\_quality\_control\_script.sh** had been executed, the folder: **09\_aligned\_reads\_quality\_control** will contain all these quality assessment reports for every **.bam** file in folder **08\_results**, just like below:



The left figure shows a single-end dataset with four ChIP samples and three control samples. The right figure shows a paired-end dataset with two ChIP samples and two control samples. The **.zip** files contain the individual components to be compiled for the report so you can ignore those. To read the reports, open the **.html** files. This file is a multitabular file in which you can evaluate how your alignment to reference genome and filtering based on mapping score affected your overall sequencing reads. Comprehensive as it is, explaining the contents in detail would take a whole new guide by itself. Therefore, in case the reports do not spell everything out enough for you, check this out: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/>.





## 16. Peak calling.

For Transcription Factors - Performed by MACS2 (default setting), GEM, HOMER (factor setting), and Genrich for transcription factor proteins of interest.

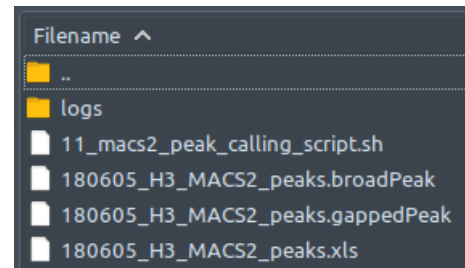
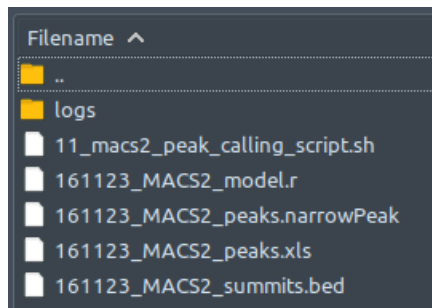
For Broad Peaks (Histone Marks) - Performed by MACS2 (broad setting), SICER2, HOMER (broad setting), and Genrich for histone modifier protein of interest.

The same track of aligned reads is scanned for potential protein-DNA binding sites. The process returns a list of enriched regions in various formats.

Modular script used: **11\_macs2\_peak\_calling\_script.sh**

- **Calls** : macs2 callpeak
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bam
- **Process** : Generate a list of called peaks
- **Output** : 11\_macs2\_peak\_calling / [setname]\_MACS2\_peaks.narrowPeak

**How does it look like?** After **11\_macs2\_peak\_calling\_script.sh** had been executed, the folder: **11\_macs2\_peak\_calling** will contain all these files, just like below:



For some reasons, MACS2 does not generate the statistical model of the dataset's background noise (.r) in paired-end mode (right figure).

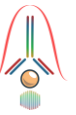
MACS2 generates three output files, which all are basically peak lists, but each of these has their own exclusive informations. In datasets with narrow peak type samples (left figure) the file that contains all necessary information relevant to analysis by ChIP-AP is the one with **.narrowPeak** extension. The **peaks.xls** file has the pileup value and peak length information, which tells us about the overall coverage of the corresponding peak region. The **summits.bed** basically are just lists of only the peak summits coordinates and read depth at the respective 1 base coordinate.

On the other hand, in datasets with broad peak type samples (right figure) the file that contains all necessary information relevant to analysis by ChIP-AP is the one with **.broadPeak** extension. The **.gappedPeak** file is basically a variant of **.broadPeak**, which is dedicated for peaks with both narrow and broad characteristics mixed in together, and thus has values that describes where and how deep are the “thick” and “thin” regions along each peak. As in the case of narrow peak type, the **peaks.xls** file has the pileup value and peak length information, which tells us about the overall coverage of the corresponding peak region. No **summits.bed** file generated in this case because of the absence of such “summit” in broad peaks.



The **.narrowPeak** file which ChIP-AP utilizes, has the following columns:

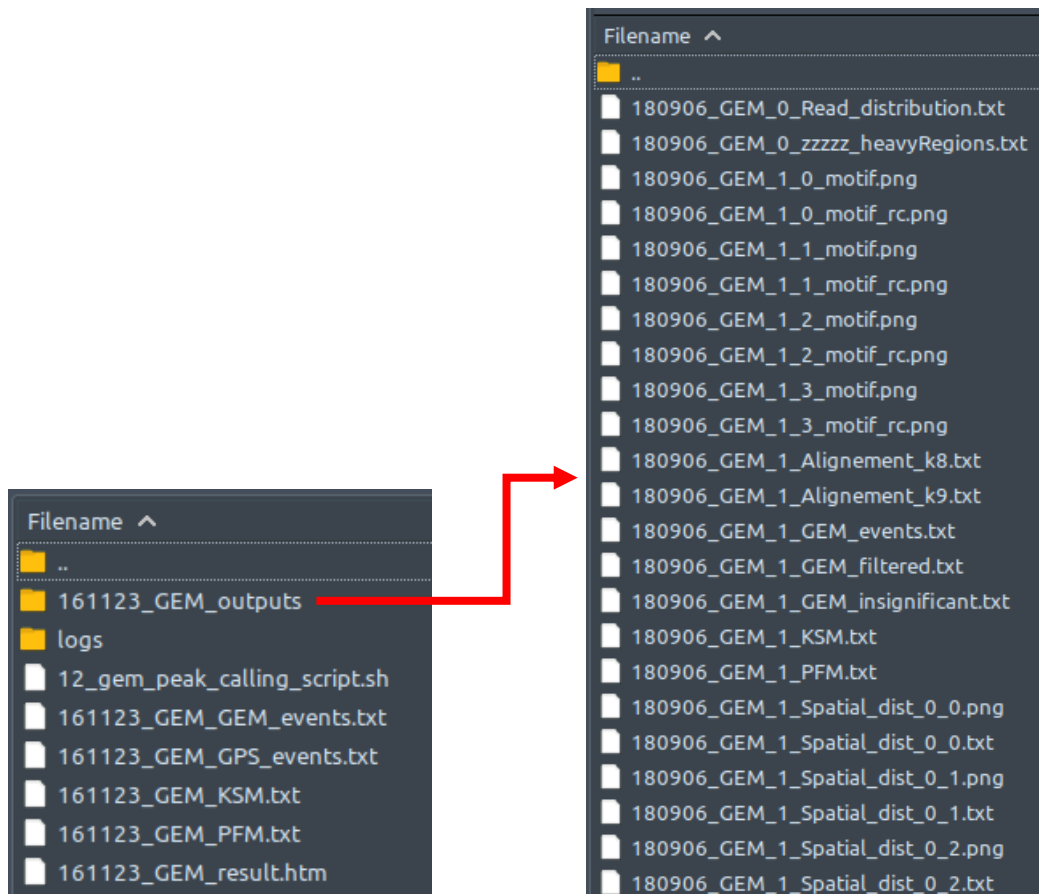
1. **chrom** - Name of the chromosome (or contig, scaffold, etc.).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chromosome is numbered 0.
3. **chromEnd** - The ending position of the feature in the chromosome or scaffold. The *chromEnd* base is not included in the display of the feature. For example, the first 100 bases of a chromosome are defined as *chromStart=0*, *chromEnd=100*, and span the bases numbered 0-99.
4. **name** - Name given to a region. Use "." if no name is assigned.
5. **score** - Indicates how dark the peak will be displayed in the browser (0-1000). If all scores were "0" when the data were submitted to the DCC, the DCC assigned scores 1-1000 based on signal value. Ideally the average signalValue per base spread is between 100-1000.
6. **strand** - +/- to denote strand or orientation (whenever applicable). Use "." if no orientation is assigned.
7. **signalValue** - Overall (usually, average) enrichment for the region.
8. **pValue** - Statistical significance (-log10). Use -1 if no pValue is assigned.
9. **qValue** - Statistical significance using false discovery rate (-log10). Use -1 if no qValue is assigned.
10. **peak** - Point-source called for this peak; 0-based offset from chromStart. Use -1 if no point-source called.



Modular script used: **12\_gem\_peak\_calling\_script.sh**

- **Calls** : gem
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bam
- **Process** : Generate a list of called peaks
- **Output** : 12\_gem\_peak\_calling / [setname]\_GEM\_GEM\_events.txt
- **Condition**: Peak type is narrow. GEM is replaced by SICER2 for broad peak type.

**How does it look like?** After **12\_gem\_peak\_calling\_script.sh** had been executed, the folder: **12\_gem\_peak\_calling** will contain all these files, just like below:



GEM outputs both the binding event files and the motif files. Because of the read distribution re-estimation, GEM outputs event prediction and read distribution files for multiple rounds. All of these are saved in folder **[setname]\_GEM\_outputs**. However, as long as ChIP-AP is concerned, we are only interested in GEM's final output files, which are saved outside the said folder. Each of these has their own exclusive information. GEM is actually a suite consisting of multiple modules performing their specific tasks. The GPS module is the one that detects peaks based on reads distribution (similar to most peak callers). The resulting peak list from solely running this GPS module can be viewed in file **[setname]\_GEM\_GPS\_events.txt**.

However, GEM is also equipped with motif enrichment analysis module that helps improve true peaks detection. This GPS peak list is then processed further and modified based on the motif enrichment analysis, resulting in the final peak list in file **[setname]\_GEM\_GEM\_events.txt**, which is the one utilized by ChIP-AP.



GEM also generates two secondary output files **[setname]\_GEM\_KSM.txt** and **[setname]\_GEM\_PFM.txt**, which are more of motif enrichment results rather than peak lists, and thus will not be discussed further here. Do check GEM documentations which link is provided at the end of this guide if you are interested. Lastly, **[setname]\_GEM\_result.htm** is a web-based comprehensive summary of all the binding events and motifs.

The **\_GEM\_events.txt** file which ChIP-AP utilizes, has the following columns:

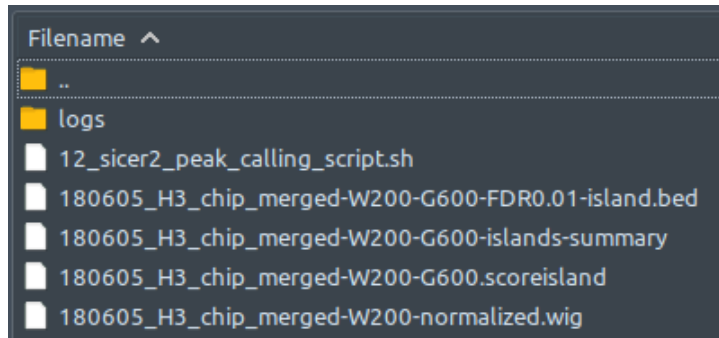
1. **Location** - The genome coordinate of this binding event
2. **IP binding strength** - The number of IP reads associated with the event
3. **Control binding strength** - the number of control reads in the corresponding region
4. **Fold** - Fold enrichment (IP/Control)
5. **Expected binding strength** - The number of IP read counts expected in the binding region given its local context (defined by parameter W2 or W3), this is used as the Lambda parameter for the Poisson test
6. **Q\_-lg10** -  $-\log_{10}(\text{q-value})$ , the q-value after multiple-testing correction, using the larger p-value of Binomial test and Poisson test
7. **P\_-lg10** -  $-\log_{10}(\text{p-value})$ , the p-value is computed from the Binomial test given the IP and Control read counts (when there are control data)
8. **P\_poiss** -  $-\log_{10}(\text{p-value})$ , the p-value is computed from the Poisson test given the IP and Expected read counts (without considering control data)
9. **IPvsEMP** - Shape deviation, the KL divergence of the IP reads from the empirical read distribution ( $\log_{10}(\text{KL})$ ), this is used to filter predicted events given the `--sd` cutoff (default=-0.40).
10. **Noise** - The fraction of the event read count estimated to be noise
11. **KmerGroup** - The group of the k-mers associated with this binding event, only the most significant k-mer is shown, the n/n values are the total number of sequence hits of the k-mer group in the positive and negative training sequences (by default total 5000 of each), respectively
12. **KG\_hgp** -  $\log_{10}(\text{hypergeometric p-value})$ , the significance of enrichment of this k-mer group in the positive vs negative training sequences (by default total 5000 of each), it is the hypergeometric p-value computed using the pos/neg hit counts and total counts
13. **Strand** - The sequence strand that contains the k-mer group match, the orientation of the motif is determined during the GEM motif discovery, '\*' represents that no k-mer is found to associated with this event



Modular script used: **12\_sicer2\_peak\_calling\_script.sh**

- **Calls** : sicer
- **Input** : 08\_results / [setname]\_chip/ctrl\_merged.bam
- **Process** : Generate a list of called peaks
- **Output** : 12\_sicer2\_peak\_calling / [setname]\_chip\_merged-W\*-G\*-islands-summary  
(\* depends on -w and -g flag arguments. Defaults are 200 and 600, respectively)
- **Condition**: Peak type is broad. SICER2 is replaced by GEM for narrow peak type.

**How does it look like?** After **12\_sicer2\_peak\_calling\_script.sh** had been executed, the folder: **12\_sicer2\_peak\_calling** will contain all these files, just like below:



SICER2 generates multiple output files as follows:

- **[setname]\_chip\_merged-W\*-G\*-.scoreisland**: delineation of significant islands controlled by E- value of 1000. It is in “chrom start end score” format.
- **[setname]\_chip\_merged-W\*-normalized.wig**: wig file that can be used to visualize the windows generated by SICER2. Read count is normalized by library size per million.
- **[setname]\_chip\_merged-W\*-G\*-islands-summary**: summary of all candidate islands with their statistical significance. It is a tab-separated-values file that has the following columns format: **chrom, start, end, ChIP\_island\_read\_count, CONTROL\_island\_read\_count, p\_value, fold\_change, FDR\_threshold**.
- **[setname]\_chip\_merged-W\*-G\*-FDR\*-island.bed**: delineation of significant islands filtered by false discovery rate (FDR). It has the following format: chrom, start, end, read-count.

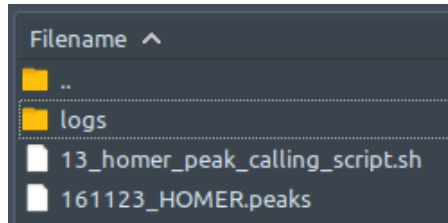
The file that contains all necessary information relevant to analysis by ChIP-AP is the one with **[setname]\_chip\_merged-W\*-G\*-islands-summary**.



Modular script used: **13\_homer\_peak\_calling\_script.sh**

- **Calls** : findPeaks
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#].bam
- **Process** : Generate a list of called peaks
- **Output** : 13\_homer\_peak\_calling / [setname]\_HOMER.peaks

**How does it look like?** After **13\_homer\_peak\_calling\_script.sh** had been executed, the folder: **13\_homer\_peak\_calling** will contain all these files, just like below:



There is no output files difference between single-end and paired-end mode.

Even though HOMER has different modes for narrow peak type (using flag argument -style factor) and broad peak type (using flag argument -style histone), the output files stay the same and so do the contents of these files, except for column 7 (see below)

**The HOMER.peaks** which CHIP-AP utilizes, has the following columns:

1. **PeakID** - Unique name for each peak
2. **chr** - Chromosome where peak is located
3. **start** - Starting position of peak
4. **end** - Ending position of peak
5. **Strand** (+/-)
6. **Normalized Tag Counts** - Number of tags found at the peak, normalized to 10 million total mapped tags (or defined by the user)
7. **Focus Ratio** - Fraction of tags found appropriately upstream and downstream of the peak center. (when sample peak type = narrow), **OR**  
**Region Size** - Length of enriched region (when sample peak type = broad)
8. **Peak score** - Position adjusted reads from initial peak region reads per position may be limited)
9. **Total Tags** - Peak depth in the ChIP sample (normalized to control)
10. **Control Tags** - Peak depth in the control sample
11. **Fold Change vs Control** - Peak depth fold change of ChIP compared to control
12. **p-value vs Control** - Statistical significance of ChIP peak compared to control
13. **Fold Change vs Local** - Peak depth fold change of ChIP compared to its surrounding regions
14. **p-value vs Local** - Statistical significance of ChIP peak compared to its surrounding regions
15. **Clonal Fold Change** - Statistical significance of ChIP peak considering the abundance of read fragment clones, or duplicates

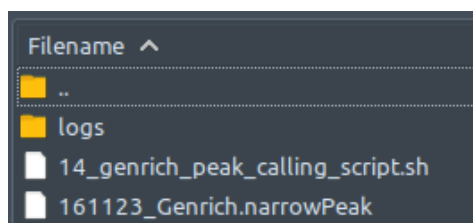




Modular script used: **14\_genrich\_peak\_calling\_script.sh**

- **Calls** : Genrich
- **Input** : 08\_results / [setname]\_[chip/ctrl]\_rep[#]\_namesorted.bam
- **Process** : Generate a list of called peaks
- **Output** : 14\_genrich\_peak\_calling / [setname]\_Genrich.narrowPeak

**How does it look like?** After **14\_genrich\_peak\_calling\_script.sh** had executed, the folder: **14\_genrich\_peak\_calling** will contain all these files, just like below:



There is no output files difference between single-end and paired-end mode.

The **.narrowPeak** file which ChIP-AP utilizes, has the following columns:

1. **chrom** - Name of the chromosome (or contig, scaffold, etc.).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chromosome is numbered 0.
3. **chromEnd** - The ending position of the feature in the chromosome or scaffold. The *chromEnd* base is not included in the display of the feature. For example, the first 100 bases of a chromosome are defined as *chromStart=0*, *chromEnd=100*, and span the bases numbered 0-99.
4. **name** - Name given to a region. Use "." if no name is assigned.
5. **score** - Indicates how dark the peak will be displayed in the browser (0-1000). If all scores were "0" when the data were submitted to the DCC, the DCC assigned scores 1-1000 based on signal value. Ideally the average signalValue per base spread is between 100-1000.
6. **strand** - +/- to denote strand or orientation (whenever applicable). Use "." if no orientation is assigned.
7. **signalValue** - Overall (usually, average) enrichment for the region.
8. **pValue** - Statistical significance (-log10). Use -1 if no pValue is assigned.
9. **qValue** - Statistical significance using false discovery rate (-log10). Use -1 if no qValue is assigned.
10. **peak** - Point-source called for this peak; 0-based offset from chromStart. Use -1 if no point-source called.

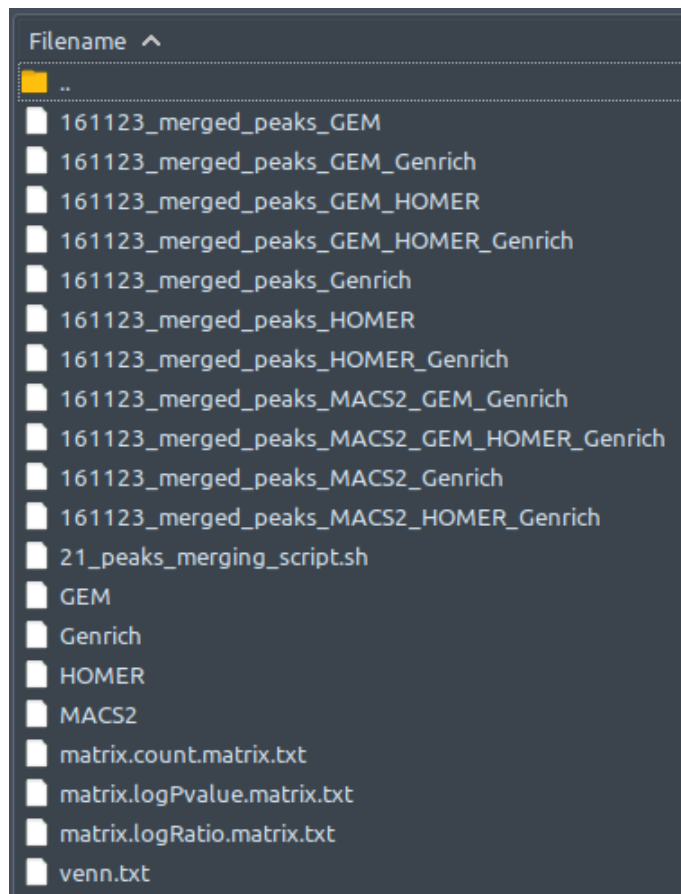


17. **Peaks merging.** Performed by a custom script and HOMER's mergePeaks. The custom script reformats necessary peak caller outputs into HOMER region list format. mergePeaks looks for overlaps between the regions in the four peak caller outputs and lists the merged regions in multiple files based on the peak caller(s) that calls them. These multiple files are then concatenated together into a single regions list file.

Modular script used: **21\_peaks\_merging\_script.sh**

- **Calls** : mergePeaks
- **Input** : 11\_macs2\_peak\_calling / [setname]\_MACS2\_peaks.narrowPeak  
12\_gem\_peak\_calling / [setname]\_GEM\_GEM\_events.txt (narrow peak only)  
12\_sicer2\_peak\_calling / [setname]-W\*-G\*-islands-summary (broad peak only)  
13\_homer\_peak\_calling / [setname]\_HOMER.peaks  
14\_genrich\_peak\_calling / [setname]\_Genrich.narrowPeak
- **Process** : Generate multiple lists of merged peak coordinates based on peak callers
- **Output** : 21\_peaks\_merging / [setname]\_merged\_peaks\*  
\* is the combination of peak callers where the listed peaks in are found in  
(e.g., [setname]\_merged\_peaks\_MACS2\_Genrich)

**How does it look like?** After **21\_peaks\_merging\_script.sh** had executed, the folder: **21\_peaks\_merging** will contain all these files, just like below:



There is no output files difference between single-end and paired-end mode. Also, you will find SICER2 instead of GEM when running broad peak type datasets.

ChIP-AP takes the input files described above, then generates the four tab-separated-values files **MACS2**, **GEM** or **SICER2**, **HOMER**, and **Genrich**, which are basically lists of peaks detected by their respective peak caller.



With **MACS2**, **GEM** or **SICER2**, **HOMER**, and **Genrich** as the input peak list, **HOMER mergePeaks** generates separate files based on overlapping peaks for each set of peaks: **21\_peaks\_merging/[setname]\_merged\_peaks\***, where \* is the combination of peak callers where the listed peaks in are found in. Files with certain peak callers combinations that contains zero peak will be non-existent in this folder, so don't be alarmed if, for example, you cannot find **[setname]\_merged\_peaks\_MACS2\_GEM** file. That simply means that there is no peak that is detected ONLY by MACS2 and GEM, just like what we can see from the example above.

The following three **matrix.txt** files below are the secondary outputs of **HOMER mergePeaks** program run. These files basically contains the statistics about the pairwise overlap of peaks between the four callers peak sets, which could provide some information for you despite not being any use for further processes down the pipeline. These explanations below are taken directly from HOMER documentation - no further explanation is given by HOMER, so we cannot give a clearer explanation:

- **matrix.logPvalue.matrix.txt**: natural log p-values for overlap using the hypergeometric distribution, positive values signify divergence
- **matrix.logRatio.matrix.txt**: natural log of the ratio of observed overlapping peaks to the expected number of overlapping peaks
- **matrix.count.matrix.txt**: raw counts of overlapping peaks

Finally, there is the file **venn.txt**. This contains the numbers needed for you to create a venn diagram depicting the peak overlaps between the four peak caller sets. Some custom scripts are available online which are able to directly take this **venn.txt** file as an input and generates a venn diagram image as a result.

Modular script used: **22\_peaks\_processing\_script.sh**

- **Operation**: cat (Bash)
- **Input** : 21\_peaks\_merging / [setname]\_merged\_peaks\*
- **Process** : Generate concatenated list of peak coordinates
- **Output** : 22\_peaks\_processing / [setname]\_all\_peaks\_concatenated.tsv

***How does it look like?** Detailed output preview of this, step 18, 19, 20, and 21 are combined below step 21*

18. **Peaks annotation**. Performed by **annotatePeaks** from HOMER package. Each region in the concatenated list is annotated based on its genomic location for the genome specified. The process returns the same list of regions, with each entry row appended with various information pertaining to the gene name, database IDs, category, and instances of motif (if HOMER known motif matrix file is provided to ChIP-AP), etc.

Modular script used: **22\_peaks\_processing\_script.sh**

- **Calls** : **annotatePeaks**
- **Input** : 22\_peaks\_processing/[setname]\_all\_peaks\_concatenated.tsv
- **Process** : Append gene annotations to the list of peak coordinates
- **Output** : 22\_peaks\_processing/[setname]\_all\_peaks\_annotated.tsv

***How does it look like?** Detailed output preview step 18, 19, 20, and 21 are combined below step 21*



19. **Fold enrichment calculations.** Performed by a custom script with the help of samtools depth and view modules. For weighted peak center fold enrichment calculation in cases of narrow peak type datasets, the custom script sends out the reformatted genomic regions as command line arguments for multi-threaded samtools depth runs. Samtools depth returns a list of read depths at each base within the region and saves them in a temporary file. The script then reads the temporary files and determine the weighted peak centers and returns the read depth values along with the base locations. The custom script sends out the weighted peak center base locations as command line arguments for multi-threaded samtools view runs. Samtools view returns the read depth values at the given base locations. The custom script then calculates the fold enrichment values, corrected based on ChIP-to-control normalization factor.

For average fold enrichment calculation in cases of broad peak type datasets, samtools view simply sums up the number of reads in the whole peak region, then calculates the fold enrichment values, corrected based on ChIP-to-control normalization factor.

In addition, the custom-made script also makes some reformatting and provides additional information necessary for downstream analysis.

Modular script used: **22\_peaks\_processing\_script.sh**

- **Calls** : fold\_change\_calculator.py
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_annotated.tsv
- **Process** : Calculate ChIP tag counts (read depth)  
Calculate weighted center fold change (narrow peak only)  
Calculate average fold change (broad peak only)  
Calculate number of peak callers overlaps  
Calculate number of user-provided (via --motif flag) motif instances found
- **Output** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv

***How does it look like?** Detailed output preview step 18, 19, 20, and 21 are combined below step 21*

20. **Irreproducibility rate (IDR) calculation.** Performed by a custom script plugged to IDR module. The IDR module compares two different peak sets and assign to each listed peak in both peak sets an irreproducibility rate (IDR) value based on that peak's capacity to be recalled by the other peak set. IDR value of each peak listed in the full (union) peak list ([setname]\_all\_peaks\_calculated.tsv) was obtained by pairing it against every individual peak caller sets, followed by calculating the pair-wise -logIDR values, then summing them all up, and finally converting it into a final IDR value. The final IDR value shows the chance of a finding (i.e., the peak) being unable to be reproduced by different peak calling algorithms. This step reprocesses the peak list [setname]\_all\_peaks\_calculated.tsv, augment the table with relevant IDR values, and re-save it under the same file name. For more details about IDR calculation method, see the IDR module documentations.

Modular script used: **22\_peaks\_processing\_script.sh**

- **Calls** : IDR\_integrator.py, IDR
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Calculate -logIDR of peaks between union peak set vs MACS2 peak set  
Calculate -logIDR of peaks between union peak set vs GEM/SICER2 peak set  
Calculate -logIDR of peaks between union peak set vs HOMER peak set  
Calculate -logIDR of peaks between union peak set vs Genrich peak set  
Calculate IDR value of peaks from the sum of all four -logIDR values
- **Output** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv  
(Ready to view, if user does not wish for gene ontology or pathway annotations)



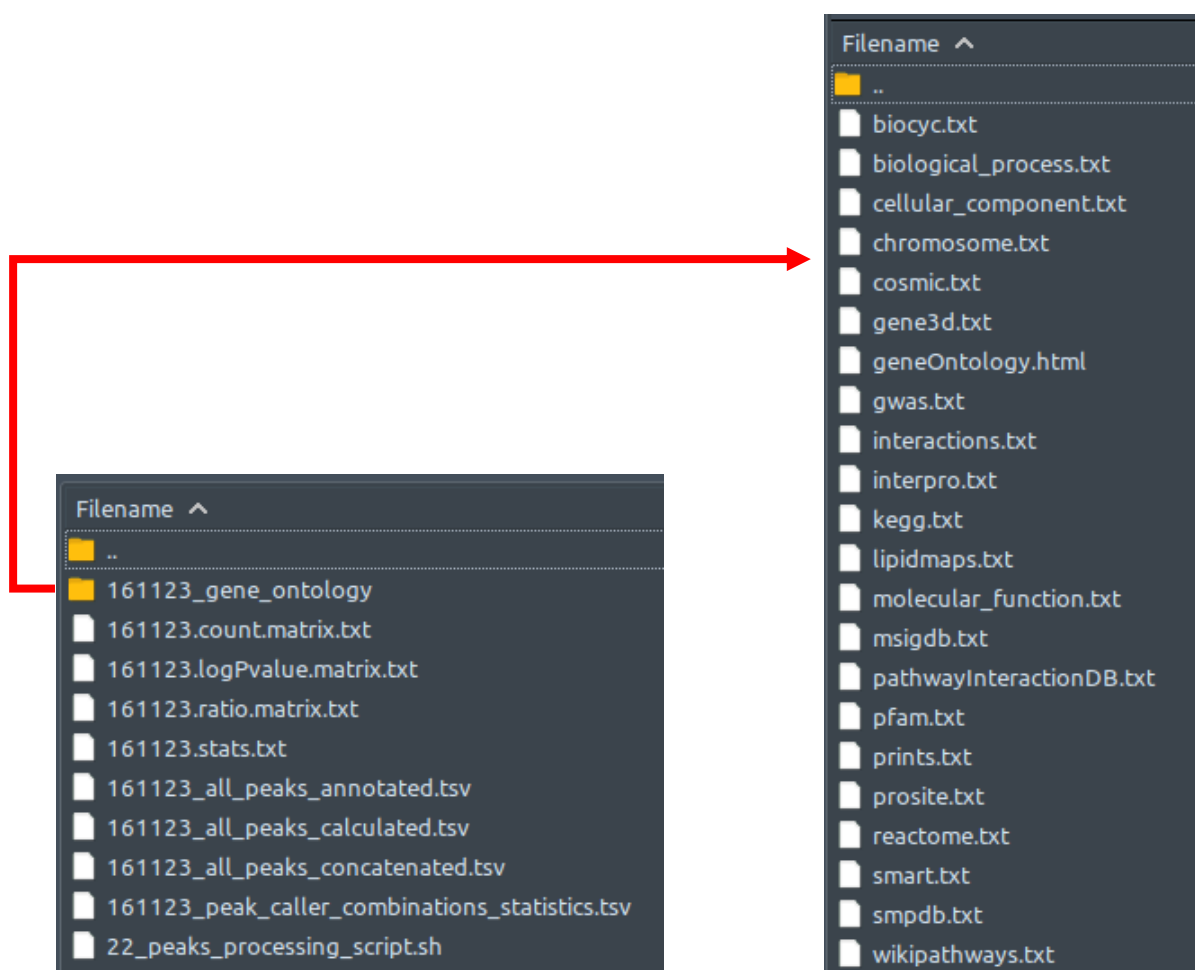
*How does it look like? Detailed output preview step 18, 19, 20, and 21 are combined below step 21*

21. **Peak statistics summary.** Performed by a custom script designed for quality assessment of called peaks. Returns a summary text file containing information pertaining to the peak read depth, peak fold enrichment, known motif hits, and positive peak hits (based on known motif presence), in each peak set along the continuum between single peak callers and the absolute consensus of all four peak callers.

Modular script used: **22\_peaks\_processing\_script.sh**

- **Calls** : peak\_caller\_stats\_calculator.py
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Generate a separate summary table of key statistics in peak callers performance
- **Output** : 22\_peaks\_processing / [setname]\_peak\_caller\_combinations\_statistics.tsv

*How does it look like?* After **22\_peaks\_processing\_script.sh** had executed, the folder: **22\_peaks\_processing** will contain all these files, just like below:



The last segment of step 17 concatenates all the peak caller combinations peak list files into one file: **[setname]\_all\_peaks\_concatenated.tsv**, followed by annotation in step 18 that generates the file: **[setname]\_all\_peaks\_annotated.tsv**, followed by fold change and IDR calculation in step 19 and 20 that generates the file: **[setname]\_all\_peaks\_calculated.tsv**. Step 21 reads the resulting file **[setname]\_all\_peaks\_calculated.tsv** and generates a statistics summary file **[setname]\_peak\_caller\_combinations\_statistics.tsv**.



While HOMER annotatePeaks are working on annotating our concatenated peak list, it also performs a gene ontology enrichment analysis which generates several tab-separated-values text files as depicted in the right figure. Each of these files contains ranked list of enriched terms coming from specific genome ontology or pathway database. As ChIP-AP will only further append **[setname]\_peak\_caller\_combinations\_statistics.tsv** with terms from certain databases (optional; activated by **--goann** and/or **--pathann** flag), not all files generated here will be used further in the pipeline. The files used are **biological\_process.txt**, **molecular\_function.txt**, **cellular\_component.txt**, **interactions.txt**, **cosmic.txt**, **kegg.txt**, **biocyc.txt**, **pathwayInteractionDB.txt**, **reactome.txt**, **smpdb.txt**, and **wikipathways.txt**.

The following three **matrix.txt** and a **stats.txt** files below are the secondary outputs of **HOMER annotatePeaks** program run. These files basically contains the statistics about the co-occurrence of motif instances (provided with **--motif** flag argument) in the peak sets, which could provide some information for you despite not being any use for further processes down the pipeline. These explanations below are taken directly from HOMER documentation - no further explanation is given by HOMER, so we cannot give a clearer explanation:

- **setname.count.matrix.txt** - number of peaks with motif co-occurrence
- **setname.ratio.matrix.txt** - ratio of observed vs. expected co-occurrence
- **setname.logPvalue.matrix.txt** - co-occurrence enrichment
- **setname.stats.txt** - table of pair-wise motif co-occurrence statistics

At this point, the results are actually ready to for your to view and analyze as they already have the essential information typically needed for ChIP-seq analysis. More details are described in the section below: **Main Pipeline Output - Final Analysis Table**. Here is a quick summary of what these are and why are they relevant to your analysis.

**[setname]\_all\_peaks\_concatenated.tsv** already has information pertaining to:

- Peak ID
- Chr
- Start
- End
- Strand
- Peak Caller Combination

So if you basically only need to know where the peaks are, and which peak caller managed to detect particular peaks, this will suffice. For example: if you want to overlap the list detected peaks with your list of genomic coordinates (e.g., of genome-wide motif instance locations, or genome-wide histone marker locations, or regions of interests obtained from different experiment, or peak list you obtained by using your favorite peak caller etc.).

**[setname]\_all\_peaks\_annotated.tsv** has these following information in addition to what is already in **[setname]\_all\_peaks\_concatenated.tsv**:

- Annotation
- Detailed Annotation
- Distance to TSS
- Nearest PromoterID
- Entrez ID
- Nearest Unigene
- Nearest Refseq





- Nearest Ensembl
- Gene Name
- Gene Alias
- Gene Description
- Gene Type
- CpG%
- GC%

At this point, the peak list is finally something biologically relevant, as each peak are now appended with the information that can be used to infer role and functionality at cellular or organism level, based on the nearest gene from the peak coordinate. As the protein used in ChIP pulldown experiments are typically transcription factor or histone modifier, a binding event in the close vicinity to a gene suggests regulation of gene expression. For instance, analyze this together with RNAseq differentially expressed genes data, then you might find which genes or which pathways your protein of interest is upregulating or downregulating.

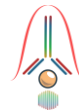
As these peaks are now also equipped by the multiple databases' ID of their nearest genes, the user can now connect this peak list with another list which entries are identified by a specific unique ID (e.g. ChIP-AP supplementary annotations relies on individual peak's Entrez ID in order to connect to HOMER genome ontology databases and subsequently add genome ontology and pathway terms into every peak in the list).

**[setname]\_all\_peaks\_calculated.tsv** has these following information in addition to what is already in **[setname]\_all\_peaks\_annotated.tsv**:

- Peak Caller Overlaps
- ChIP Tag Count
- Control Tag Count
- Fold Change
- Number of Motifs

At this point, you have more power to evaluate and select the peaks you want in your final set. In this file you now have the actual read depth of your peaks, and also the fold change value where you can see the enrichment of reads at the potential binding site compared to the control sample with no pulldown. Along with that, each peak also have the number of known DNA binding motif found in the sample. Note that this value will only appear if you provided the **.motif** file using the **--motif** flag argument. These values provided here are the most basic properties of peaks pertaining to their confidence, and are the most standard ways of filtering and ranking of peaks in the list. These values are provided here as alternative ways to apply some thresholds in order to select your peak set for further analysis, in addition to the more powerful, main method provided by this pipeline: multiple peak caller overlaps.

To accomodate, at this point, the peak list now has the number of peak caller overlaps, which is simply the number of peak callers that detected this peak. Although user can already filter in or out their peak list based on the peak caller names provided by the column "Peak Caller Combination" in the file **[setname]\_all\_peaks\_concatenated.tsv**, this value is here to give user a more convenient way of filtering their peaks based on how many peak callers "agree" with a particular detected peak, regardless of which peak callers detected it.



22. **(Optional) Downstream analysis: Gene ontology enrichment.** Each peak in the concatenated list is appended with all the gene ontology terms associated with its gene annotation. The gene ontology terms are derived from biological processes, molecular functions, and cellular compartments databases. This enables list filtering based on the gene ontology terms of the study's interest

Modular script used: **23\_go\_annotation\_script.sh**

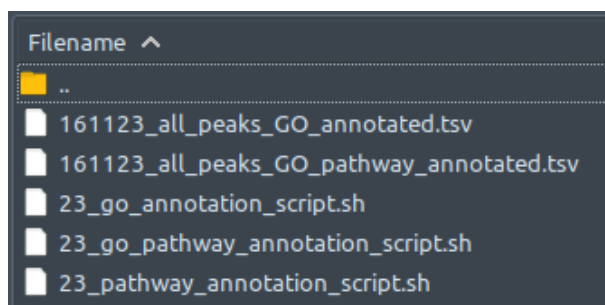
- **Calls** : go\_annotator.py
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Append related gene ontology terms to the list of peaks
- **Output** : 23\_supplementary\_annotations / [setname]\_all\_peaks\_go\_annotated.tsv

23. **(Optional) Downstream analysis: Pathway enrichment.** Each peak in the concatenated list is appended with all the related biological pathways associated with its gene annotation. The biological pathway terms are derived from KEGG, SMPDB, Biocyc, Reactome, Wikipathways, and pathwayInteractionDB databases. This enables list filtering based on the biological pathways of the study's interest. Additionally, this analysis also adds other terms pertaining to known interactions with common proteins and known gene mutations found in malignant cases, derived from common protein interaction and COSMIC databases, respectively.

Modular script used: **23\_pathway\_annotation\_script.sh**

- **Calls** : pathway\_annotator.py
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Append known pathways and interactions to the list of peaks
- **Output** : 23\_supplementary\_annotations / [setname]\_all\_peaks\_pathway\_annotated.tsv

**How does it look like?** After **23\_go\_annotation\_script.sh** had executed, the folder: **23\_supplementary\_annotations** will contain all these files, just like below:



There is no output files difference between single-end and paired-end mode.

If **--goann** flag is used during ChIP-AP call, **23\_go\_annotation\_script.sh** will be executed, and generates **[setname]\_all\_peaks\_go\_annotated.tsv** that contains the following gene ontology terms based on each peak's nearest gene:

- Biological Process
- Molecular Function
- Cellular Component



If **--pathann** flag is used during ChIP-AP call, **23\_pathway\_annotation\_script.sh** will be executed, and generates **[setname]\_all\_peaks\_pathway\_annotated.tsv** that contains the following known pathways terms based on each peak's nearest gene:

- Interaction with Common Protein
- Somatic Mutations (COSMIC)
- Pathway (KEGG)
- Pathway (BIOCYC)
- Pathway (pathwayInteractionDB)

If both **--goann** flag and **--pathann** flag are used during ChIP-AP call, both **23\_go\_annotation\_script.sh** and **23\_pathway\_annotation\_script.sh** will be executed, and generates **[setname]\_all\_peaks\_go\_pathway\_annotated.tsv** that contains all the information that are gained by running the two scripts one after another:

- Biological Process
- Molecular Function
- Cellular Component
- Interaction with Common Protein
- Somatic Mutations (COSMIC)
- Pathway (KEGG)
- Pathway (BIOCYC)
- Pathway (pathwayInteractionDB)

These columns contains **ALL** the related terms from their respective gene ontology or pathway databases, which are comma-separated between terms. That said, depending on how developed the databases are, and how much is known about the gene, the amount of information can range between none to overwhelming.

The information is very useful for filtering the peak list based on the protein of interest's potential roles in specific biological activities or pathways, based on the interest of the study. This can also be very handy, because contrary to the conventional way of only looking at the gene ontology enrichment analysis result and manually checking back-and-forth if specific genes of interest are actually related to specific terms of interest, we have it already linked to each peak in the list.

However, with respect to users who do not need such information and probably want their list to be much less verbose and smaller in size, this step is completely optional. The output files will not be generated, to save processing time should the user choose to omit this step (by not using the respective flags or not ticking the boxes in GUI). The scripts will still be generated by ChIP-AP, though, just not executed. You can simply run the script should you change your mind and decide to have these supplementary annotations.

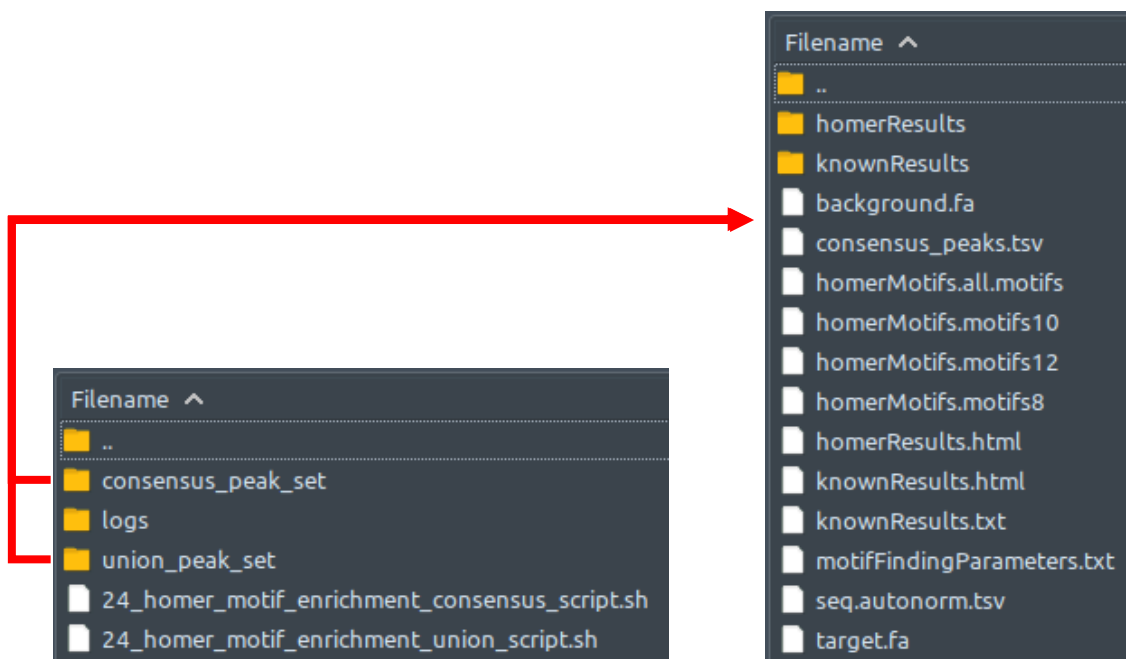


24. **(Optional) Downstream analysis: Motif enrichment analysis with HOMER.** Genomic sequences are extracted based on the coordinates of the peaks in consensus (four peak callers overlap), union (all called peaks), or both peak lists. HOMER performs analysis to identify specific DNA sequence motifs to which the experimented protein(s) have binding affinity towards. For the sake of processing speed, HOMER utilizes cumulative binomial distribution to calculate motif enrichment by default. However, by utilizing ChIP-AP custom setting table, user may choose to utilize cumulative hypergeometric distribution, which describes motif enrichment problem more accurately. Besides the typically performed calculations for de novo motifs discovery, HOMER also calculates the enrichment scores of the known motifs in HOMER motifs database. This optional downstream analysis option is only available for datasets with narrow (transcription factor) peaks.

Modular script used: **24\_homer\_motif\_enrichment\_consensus\_script.sh**  
**24\_homer\_motif\_enrichment\_union\_script.sh**

- **Calls** : findMotifsGenome.pl
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Performs motif enrichment analysis
- **Output** : 24\_homer\_motif\_enrichment / ... / homerResults.html  
24\_homer\_motif\_enrichment / ... / knownResults.html

**How does it look like?** After **24\_homer\_motif\_enrichment\_consensus\_script.sh** or **24\_homer\_motif\_enrichment\_union\_script.sh** (or both scripts) had executed, the folder: **24\_homer\_motif\_enrichment** will contain these files, just like below:



If **--homer\_motif consensus** flag and argument are used during ChIP-AP call, **24\_homer\_motif\_enrichment\_consensus\_script.sh** will be executed, and generates **HOMER findMotifsGenome.pl** output folders and files depicted by the image to the right, inside folder **consensus\_peak\_set**. If **--homer\_motif union** flag and argument are used during ChIP-AP call, **24\_homer\_motif\_enrichment\_union\_script.sh** will be executed, and generates similar output inside folder **union\_peak\_set**. If **--homer\_motif both** flag and argument are used during ChIP-AP call, both aforementioned scripts will be executed separately, generating output files and folders inside their respective folders.

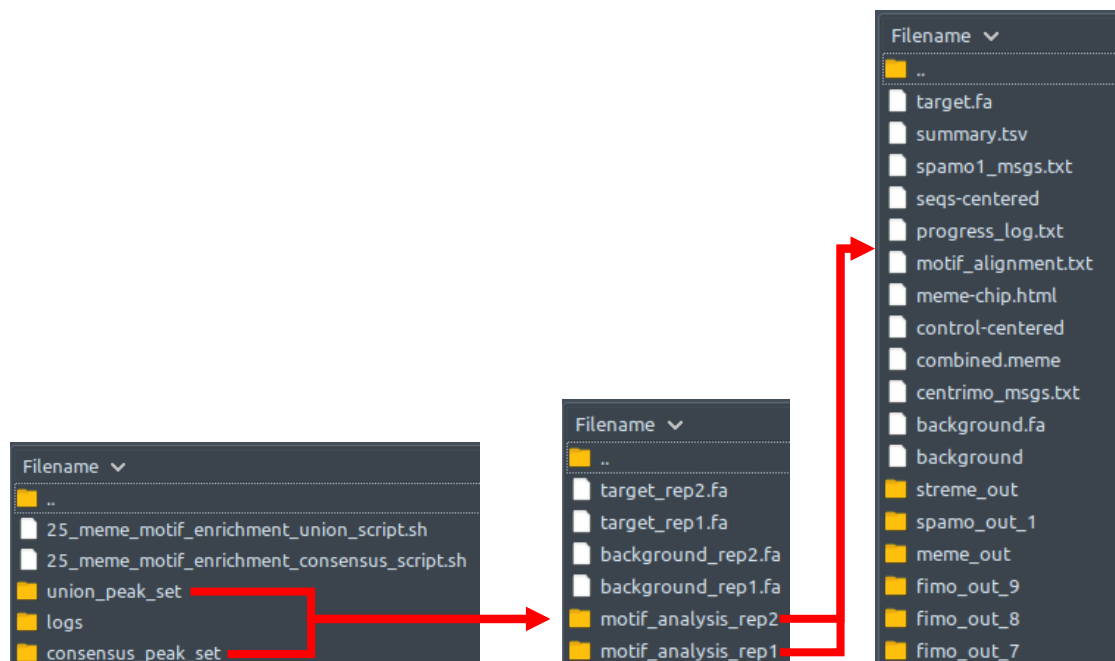


25. **(Optional) Downstream analysis: Motif enrichment analysis with MEME.** Genomic sequences are extracted based on the coordinates of the peaks in consensus (four peak callers overlap), union (all called peaks), or both peak lists. With or without control sequences extracted from random genomic sequences, MEME performs analysis to identify specific DNA sequence motifs to which the experimented protein(s) have binding affinity towards. By utilizing separate dedicated modules included in MEME suite, MEME-ChIP performs de novo motif discovery, motif enrichment analysis, motif location analysis and motif clustering in one go, providing a comprehensive picture of the DNA motifs that are enriched in the extracted sequences. MEME-ChIP performs two complementary types of de novo motif discovery: weight matrix–based discovery for high accuracy, and word-based discovery for high sensitivity. This optional downstream analysis option is only available for datasets with narrow (transcription factor) peaks.

Modular script used: **25\_meme\_motif\_enrichment\_consensus\_script.sh**  
**25\_meme\_motif\_enrichment\_union\_script.sh**

- **Calls** : meme-chip
- **Input** : 22\_peaks\_processing / [setname]\_all\_peaks\_calculated.tsv
- **Process** : Performs motif enrichment analysis
- **Output** : 25\_meme\_motif\_enrichment / ... / meme-chip.html

**How does it look like?** After **25\_meme\_motif\_enrichment\_consensus\_script.sh** or **25\_meme\_motif\_enrichment\_union\_script.sh** (or both scripts) had executed, the folder: **25\_meme\_motif\_enrichment** will contain these files, just like below:



If **--meme\_motif consensus** flag and argument are used during ChIP-AP call, **25\_meme\_motif\_enrichment\_consensus\_script.sh** will be executed, and generates **meme-chip** output folders and files depicted by the image in the middle, inside folder **consensus\_peak\_set**. If **--meme\_motif union** flag and argument are used during ChIP-AP call, **25\_meme\_motif\_enrichment\_union\_script.sh** will be executed, and generates similar output inside folder **union\_peak\_set**. If **--meme\_motif both** flag and argument are used during ChIP-AP call, both aforementioned scripts will be executed separately, generating output files and folders inside their respective folders.



When ChIP and control aligned reads (.bam) have the same number of replicates, ChIP-AP gives the option for merged (with **--fcmerge** flag) or pair-wise (without **--fcmerge** flag) fold change calculations. In pair-wise mode, peaks every ChIP vs control replicate have different weighted peak center coordinate, which directly affects the actual target and background sequences for meme-chip to perform enrichment analysis on. Therefore, ChIP-AP recognizes and processes multiple replicates separately (based on each weighted peak center coordinate), generating respective results for each replicate, stored in separate folders. On the contrary, whenever user choose to use **--fcmerge** flag, or when the number of ChIP and control samples are not the same, ChIP-AP will be forced to perform merged fold change calculation. In this situation every peak will have one weighted peak center coordinate and thus there will only be one replicate of meme-chip motif enrichment analysis results inside one single folder.





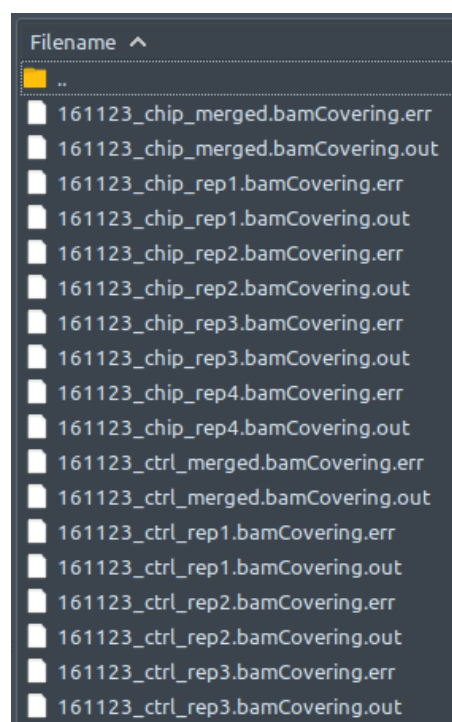
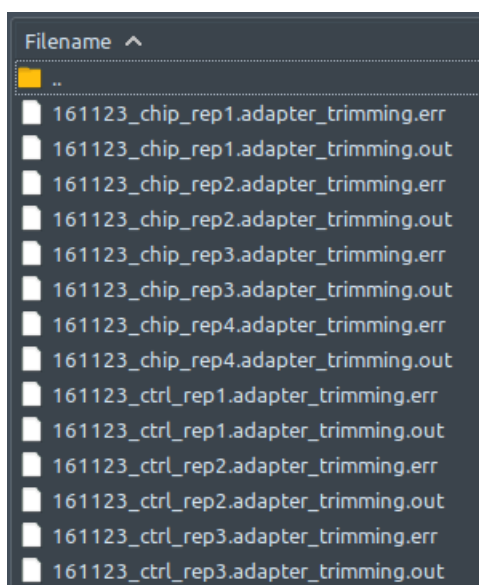
## Log Files

For most processes in every modular scripts, log files are recorded and saved in folder: /logs under their respective directories. There are two types of log files:

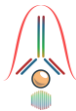
- Log files with **.out** extension: captures whatever the program writes out through channel 1>, a.k.a. the standard output
- Log files with **.err** extension: captures whatever the program writes out through channel 2>, a.k.a. the standard error

Sometimes, unlike what the file or channel name suggests, you might find errors reported in the **.out** files or something like normal program run progress report written in the **.err** files. That is just the way it is. Some programs do not follow the standard output / standard error convention, that's why. So if your pipeline crashed at a certain process and the **.err** log file does not show anything wrong, the error message might be in the **.out** file instead!

The example contents inside folder /logs can be seen below:



All these outputs are recorded for your convenience in troubleshooting and error reporting. Do open and read these files to see what's happening with your program. After you spot the potential problem, you can either post the logged error message in our GitHub - issues (<https://github.com/JSuryatenggara/ChIP-AP/issues>) , or go hit google search if you think the problem is simple and quick enough to figure out yourself.



## Main Pipeline Output

### Final Analysis Table (including supplementary annotations)

The table below shows the contents of *[filename]\_all\_peaks\_go\_pathway\_annotated.tsv*. Smaller sized and less verbose variants of this table are saved in the output folder with suffixes: concatenated, annotated, and calculated (see Source in the table below)

Column #	Peak Attribute	Source
<b>Column 1 (A)</b>	Peak ID (unique peak ID)	Pipeline script: 22_peaks_processing_script.sh Called program: cat (Bash) Output file: [filename]_all_peaks_concatenated.tsv Output folder: 22_peaks_processing
<b>Column 2 (B)</b>	Chr (chromosome)	
<b>Column 3 (C)</b>	Start (peak start coordinate)	
<b>Column 4 (D)</b>	End (peak end coordinate)	
<b>Column 5 (E)</b>	Strand (strand on which peak is found)	
<b>Column 6 (F)</b>	Peak Caller Combination	
<b>Column 7 (G)</b>	Peak Caller Overlaps	Pipeline script: 22_peaks_processing_script.sh Called script: fold_change_calculator.py IDR_integrator.py Called program: samtools depth samtools view IDR Output file: [filename]_all_peaks_calculated.tsv Output folder: 22_peaks_processing
<b>Column 8 (H)</b>	ChIP Tag Count	
<b>Column 9 (I)</b>	Control Tag Count	
<b>Column 10 (J)</b>	Fold Change	
<b>Column 11 (K)</b>	Peak Center	
<b>Column 12 (L)</b>	Number of Motifs	
<b>Column 13 (M)</b>	negLog10_IDR	
<b>Column 14 (N)</b>	IDR	
<b>Column 15 (O)</b>	Annotation	Pipeline script: 22_peaks_processing_script.sh Called program: HOMER annotatePeaks Output file: [filename]_all_peaks_annotated.tsv Output folder: 22_peaks_processing
<b>Column 16 (P)</b>	Detailed Annotation	
<b>Column 17 (Q)</b>	Distance to TSS	
<b>Column 18 (R)</b>	Nearest PromoterID	
<b>Column 19 (S)</b>	Entrez ID	
<b>Column 20 (T)</b>	Nearest Unigene	
<b>Column 21 (U)</b>	Nearest Refseq	
<b>Column 22 (V)</b>	Nearest Ensembl	
<b>Column 23 (W)</b>	Gene Name	
<b>Column 24 (X)</b>	Gene Alias	
<b>Column 25 (Y)</b>	Gene Description	
<b>Column 26 (Z)</b>	Gene Type	
<b>Column 27 (AA)</b>	CpG%	
<b>Column 28 (AB)</b>	GC%	
<b>Column 29 (AC)</b>	Biological Process	Pipeline script: 23_go_annotation_script.sh Called script: GO_annotator.py Output file: [filename]_all_peaks_go_annotated.tsv Output folder: 23_supplementary_annotations
<b>Column 30 (AD)</b>	Molecular Function	
<b>Column 31 (AE)</b>	Cellular Component	Pipeline script: 23_pathway_annotation_script.sh Called script: pathway_annotator.py Output file: [filename]_all_peaks_pathway_annotated.tsv Output folder: 23_supplementary_annotations
<b>Column 32 (AF)</b>	Interaction with Common Protein	
<b>Column 33 (AG)</b>	Somatic Mutations (COSMIC)	
<b>Column 34 (AH)</b>	Pathway (KEGG)	
<b>Column 35 (AI)</b>	Pathway (BIOCYC)	
<b>Column 36 (AJ)</b>	Pathway (pathwayInteractionDB)	
<b>Column 37 (AK)</b>	Pathway (REACTOME)	
<b>Column 38 (AL)</b>	Pathway (SMPDB)	
<b>Column 39 (AM)</b>	Pathway (Wikipathways)	



Peak Attribute	Description
Peak ID	Given unique peak ID
Chr	Chromosome where the peak is located
Start	Starting coordinate of the peak region
End	End coordinate of the peak region
Strand	DNA strand (positive/negative) where the peak is located
Peak Caller Combination	Name of peak callers that detected this peak
Peak Caller Overlaps	Number of peak callers that detected this peak
ChIP Tag Count	ChIP Read depth at weighted peak center (narrow peak) OR average read depth of the whole peak region (broad peak)
Control Tag Count	Control read depth at weighted peak center (narrow peak) OR average read depth of the whole peak region (broad peak)
Fold Change	ChIP vs control fold change at weighted peak center (narrow peak) OR ChIP vs control fold change of the whole peak region (broad peak)
Peak Center	Genomic coordinate of the weighted peak center (narrow peak) OR genomic coordinate of the of peak region midpoint (broad peak)
Number of Motifs	Number of motif instances found in the peak region
negLog10_IDR	Sum of $-\log_{10}$ IDR obtained from peak IDR calculation between the union and each individual peak caller set
IDR	Peak irreproducibility rate converted from negLog10_IDR column value
Annotation	Type of annotated region (Exon, Intron, Promoter-TSS)
Detailed Annotation	Detailed, longer version of Annotation in previous column
Distance to TSS	Distance from the peak to the nearest annotated Transcription Start Site (TSS) in basepairs
Nearest PromoterID	PromoterID of the nearest annotated promoter region
Entrez ID	Entrez ID of the nearest annotated gene
Nearest Unigene	Unigene ID of the nearest annotated gene
Nearest Refseq	Refseq ID of the nearest annotated gene
Nearest Ensembl	Ensembl ID of the nearest annotated gene
Gene Name	Name of the nearest annotated gene (abbreviated)
Gene Alias	Alternate names of the nearest annotated gene
Gene Description	Name of the nearest annotated gene (full)
Gene Type	e.g protein-coding / non-coding / pseudogene / etc
CpG%	The proportion of known methylation spots within this peak region
GC%	The proportion of GC contents within this peak region
Biological Process	Terms from biological_process.txt that are related to the nearest gene
Molecular Function	Terms from molecular_function.txt that are related to the nearest gene
Cellular Component	Terms from cellular_component.txt that are related to the nearest gene
Interaction with Common Protein	Known interactions from interaction.txt with the nearest gene
Somatic Mutations (COSMIC)	Known cancer cases where mutation of the nearest gene are found
Pathway (KEGG)	Known pathways according to the KEGG database
Pathway (BIOCYC)	Known pathways according to the Biocyc database
Pathway (pathwayInteractionDB)	Known pathways from the pathwayInteractionDB database
Pathway (REACTOME)	Known pathways according to the REACTOME database
Pathway (SMPDB)	Known pathways according to the SMPDB database
Pathway (Wikipathways)	Known pathways according to the wikipathways database



## Motif enrichment analysis results (by HOMER)

All the results are compiled and can be viewed by opening the file homerResults.html in an HTML file viewer such as your internet browser. This file gives you a formatted, organized view of the enriched de novo motifs and all the relevant information, as can be seen below. Additionally, more details can be accessed by simply clicking on the links in the table.

### Homer *de novo* Motif Results

[Known Motif Enrichment Results](#)

[Gene Ontology Enrichment Results](#)

If Homer is having trouble matching a motif to a known motif, try copy/pasting the matrix file into [STAMP](#)

More information on motif finding results: [HOMER](#) | [Description of Results](#) | [Tips](#)

Total target sequences = 37301

Total background sequences = 35962

\* - possible false positive

Rank	Motif	P-value	log P-value	% of Targets	% of Background	STD(Bg STD)	Best Match/Details	Motif File
1		1e-12661	-2.915e+04	70.91%	15.19%	40.5bp (65.1bp)	Foxa2(Forkhead)/Liver-Foxa2-ChIP-Seq/Homer <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
2		1e-578	-1.332e+03	27.14%	16.52%	54.0bp (65.5bp)	NF1-halfsite(CTF)/LNCaP-NF1-ChIP-Seq/Homer <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
3		1e-384	-8.860e+02	17.77%	10.53%	53.9bp (62.1bp)	Unknown/Homeobox/Limb-p300-ChIP-Seq/Homer <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
4		1e-164	-3.783e+02	3.17%	1.28%	52.2bp (62.9bp)	PH0048.1_Hoxa13 <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
5		1e-151	-3.485e+02	3.38%	1.47%	50.2bp (65.4bp)	NF-E2(bZIP)/K562-NFE2-ChIP-Seq/Homer <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
6		1e-107	-2.485e+02	1.21%	0.35%	56.3bp (69.7bp)	CTCF(Zf)/CD4+-CTCF-ChIP-Seq/Homer <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>
7		1e-72	-1.671e+02	2.10%	1.02%	55.1bp (58.5bp)	MA0029.1_Evi1 <a href="#">More Information</a>   <a href="#">Similar Motifs Found</a>	<a href="#">motif file (matrix)</a>

We can see these information listed below from the table:

Rank	Motif Rank
Motif	Motif position weight matrix logo
P-value	Final enrichment p-value
log P-value	Log of p-value
% of Targets	Number of target sequences with motif/ total targets
% of Background	Number of background sequences with motif/ total background
STD(Bg STD)	Standard deviation of position in target and background sequences
Best Match/Details	Best match of de novo motif to motif database

In addition to de novo motif enrichment, homer also performs motif enrichment analysis on the known binding motifs readily available within their database repertoire. The results for this analysis can be viewed in a similar way by opening the file knownResults.html, which contains similar information as its de novo counterpart.

More detailed information is available in <http://homer.ucsd.edu/homer/ngs/peakMotifs.html>



## Motif enrichment analysis results (by MEME)

All the results are compiled and can be viewed by opening the file `meme-chip.html` in an HTML file viewer such as your internet browser. This file gives you a formatted, organized view of the enriched de novo motifs and all the relevant information. Additionally, more details can be accessed by simply clicking on the links in the table.

In datasets where there is an equal number of multiple-replicated ChIP and control samples, ChIP-AP will have MEME perform a pair-wise motif enrichment analysis. Therefore, in that case, there will be multiple replicates of motif enrichment results, each one looking like below.



While the file above gives a more graphical representation of the results, `meme-chip` also generates another file: `summary.tsv`, which contains the same information (see below) repackaged in table format suitable for subsequent processing, if needed.

MOTIF_INDEX	The index of the motif in the "Motifs in MEME text format" file ('combined.meme') output by MEME-ChIP.
MOTIF_SOURCE	The name of the program that found the de novo motif, or the name of the motif file containing the known motif.
MOTIF_ID	The name of the motif, which is unique in the motif database file.
ALT_ID	An alternate name for the motif, which may be provided in the motif database file.
CONSENSUS	The ID of the de novo motif, or a consensus sequence computed from the letter frequencies in the known motif (as described below).
WIDTH	The width of the motif.
SITES	The number of sites reported by the de novo program, or the number of "Total Matches" reported by CentriMo.
E-VALUE	The statistical significance of the motif.
E-VALUE_SOURCE	The program that reported the E-value.
MOST_SIMILAR_MOTIF	The known motif most similar to this motif according to Tomtom.
URL	A link to a description of the most similar motif, or to the known motif.

More detailed information is available in <https://meme-suite.org/meme/doc/meme-chip.html>



## Miscellaneous Pipeline Outputs

### Multiple peak callers statistics summary

The table below shows the contents of `[filename]_peak_caller_combinations_statistics.tsv`.

Column 1 (A)	Peak Callers Combination
Column 2 (B)	Exclusive Peak Count
Column 3 (C)	Exclusive Positive Peak Count
Column 4 (D)	Exclusive Motif Count
Column 5 (E)	Exclusive Positive Peak Hit Rate
Column 6 (F)	Exclusive Motif Hit Rate
Column 7 (G)	Exclusive ChIP Peak Read Depth
Column 8 (H)	Exclusive ChIP Peak Fold Change
Column 9 (I)	Inclusive Peak Count
Column 10 (J)	Inclusive Positive Peak Count
Column 11 (K)	Inclusive Motif Count
Column 12 (L)	Inclusive Positive Peak Hit Rate
Column 13 (M)	Inclusive Motif Hit Rate
Column 14 (N)	Inclusive ChIP Peak Read Depth
Column 15 (O)	Inclusive ChIP Peak Fold Change

### Description:

Peak Callers Combination	Name of peak callers that detected this peaks subset
Peak Count	Number of peaks
Positive Peak Count	Number of peaks containing known binding motif (--motif)
Motif Count	Total number of binding motif instances
Positive Peak Hit Rate	Positive Peak Count divided by Peak Count
Motif Hit Rate	Motif Count divided by Peak Count
ChIP Peak Read Depth	ChIP Read depth at weighted peak center (narrow peak) OR average read depth of the whole peak region (broad peak)
ChIP Peak Fold Change	ChIP vs control fold change in read depth

- Exclusive: Only counts for a specific peak caller combination (e.g., Exclusive peak count of MACS2 only counts for peaks that is exclusively called by MACS2 alone).
- Inclusive: Counts for other peak caller combinations containing the same peak callers (e.g., Inclusive peak count of MACS2|GEM also counts for all other peaks in MACS2|GEM|HOMER, MACS2|GEM|Genrich, and MACS2|GEM|HOMER|Genrich).





## Pipeline Run Info

This file summarizes the assignment of the files (IP sample or control, read 1 or 2; replicate number) and the file name conversion for every unaligned or aligned sequencing reads to be processed. Each line tells the user what the original files have been renamed into. Check this file if you suspect the order of samples were incorrectly entered (ie swapped chip with control)

```
Chromatin IP dataset replicate 1, 1st read : Original filename = a.fastq --> New filename = setname_chip_rep1_R1.fq.gz
Chromatin IP dataset replicate 2, 1st read : Original filename = b.fastq --> New filename = setname_chip_rep2_R1.fq.gz
Chromatin IP dataset replicate 1, 2nd read : Original filename = c.fastq --> New filename = setname_chip_rep1_R2.fq.gz
Chromatin IP dataset replicate 2, 2nd read : Original filename = d.fastq --> New filename = setname_chip_rep2_R2.fq.gz
Control dataset replicate 1, 1st read : Original filename = e.fastq --> New filename = setname_ctrl_rep1_R1.fq.gz
Control dataset replicate 2, 1st read : Original filename = f.fastq --> New filename = setname_ctrl_rep2_R1.fq.gz
Control dataset replicate 1, 2nd read : Original filename = g.fastq --> New filename = setname_ctrl_rep1_R2.fq.gz
Control dataset replicate 2, 2nd read : Original filename = h.fastq --> New filename = setname_ctrl_rep2_R2.fq.gz
```

## Pipeline Run Command

Contains the input command line that was used to call the pipeline in a text file: *[filename]\_command\_line.txt* in the output save folder. This is useful for documentation of the run, and for re-running of the pipeline after a run failure or some tweaking if need be.

```
[chipap directory]/chipap.py --mode paired --ref [genome_build] --genome [path_to_computed_genome_folders]
--output [full_path_to_output_save_folder] --setname [dataset name] --sample_table [path_to_sample_table_file]
--custom_setting_table [path_to_setting_table_file].tsv --motif [path_to_known_motif_file]
--fcmerge --goann --pathann --deltemp --thread [#_of_threads_to_use] --run
```

## Sample Table

Contains the full path of each input ChIP and control sample in the pipeline run in a tab-separated value file: *[filename]\_sample\_table.tsv* in the output save folder in ChIP-AP sample table format. This is useful for documentation of the run, and for re-running of the pipeline after a run failure or some tweaking if need be. Below is an example of sample table file content (header included), given paired-end samples with two ChIP and two control replicates.

chip_read_1	chip_read_2	ctrl_read_1	ctrl_read_2
... /a.fastq	... /c.fastq	... /e.fastq	... /g.fastq
... /b.fastq	... /d.fastq	... /f.fastq	... /h.fastq

If your sample is single-ended, then the sample table can simply be formatted as follows.

chip_read_1	ctrl_read_1
... /a.fastq	... /e.fastq
... /b.fastq	... /f.fastq



## Setting Table & Default Parameters

A *cornerstone* of ChIP-AP's functionality is the settings table. ChIP-AP, with the raw fq files and the settings table, is able to reproduce any analysis (provided the same program version numbers are used). The provision of the settings table ensures reproducibility of any analysis with minimal effort and bypasses the usually sparse and significantly under-detailed methods sections of publications. Science is supposed to be reproducible, yet bioinformatics analysis are typically black-boxes which are irreproducible. This 1 file, changes that!

The structure of the settings table is simple. It is a 2 column tab-separated value file with the names of the programs on the 1<sup>st</sup> column, and the necessary flags required or changed in the 2<sup>nd</sup> column. If making your own custom table, then the 1<sup>st</sup> column below must be copied as-is and not changed. These 2 columns together, list the flags and argument values for each program used in the pipeline.

When ChIP-AP is run, a copy of the used settings table is saved as a tab-separated value file: *[filename]\_setting\_table.tsv* in the output save folder. If you have a custom settings table made and provided it as input, then ChIP-AP will make a 2<sup>nd</sup> copy of this table in the same output save folder. This decision is made as it is useful documentation of the run performed. This file is also useful for re-running of the pipeline after run failure or some tweaking if necessary. If submitting an issue request on Github, you ***must*** provide us your settings table used as well as all other requested information. See Github for details regarding this.

*We consider the dissemination of the information of this file as vital and essential along with results obtained. The table can be included as a supplemental table in a manuscript or can be included as a processed data file when submitting data to GEO – either way, the information of this file must be presented when publishing data.*

Below is an example of setting table file in its default-setting state:

fastqc1	
clumpify	dedupe spany addcount qout=33 fixjunk
bbduk	ktrim=l hdist=2
trimmomatic	LEADING:20 SLIDINGWINDOW:4:20 TRAILING:20 MINLEN:20
fastqc2	
bwa_mem	
samtools_view	-q 20
plotfingerprint	
fastqc3	
macs2_callpeak	
gem	-Xmx10G --k_min 8 --k_max 12
sicer2	
homer_findPeaks	
genrich	--adjustp -v
homer_mergePeaks	
homer_annotatePeaks	
fold_change_calculator	--normfactor uniquely_mapped
homer_findMotifsGenome	-size given -mask
meme_chip	-meme-nmotifs 25



## Interpreting ChIP-AP Output

Ok so ChIP-AP does report a fair amount of stuff. If you ran it locally you have a swath of folders and you have nooooo clue what to look for and its all confusing. We get that. The reality though its very simple to know what to look for to know your experimental run worked and in this section were going to walk you through that!

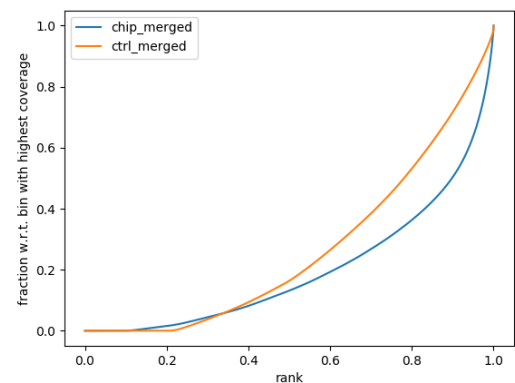
### Did my analysis work?

There are a couple of things to look for to answer this question. 1, the fingerprint plot and 2, the venn diagram of the merged peaks. Lets begin...

#### 1 – The fingerprint plot

The fingerprint plot tells us how well the enrichment of your samples worked. It is generated by the function from the deeptools package and is generated after the alignment files. As such, the plots are found in the “08\_results” folder and are labelled “fingerprint\_XXXXX.png/svg.” The PNG files allow you to view them in any image viewer, the SVG files are for opening in Adobe Illustrator or Inkscape to make HQ publication figures later if you need.

To interpret the fingerprint plot, (more information can be found on the deeptools documentation site), but put simply (image to the right), the input control should be a diagonal line as close as possible toward the 1:1 diagonal. Your ChIP sample should have a bend/kink towards the bottom right corner. The greater the separation between the input and the chip sample, the greater the enrichment you will see in the final result (ie lots of peaks). If the lines are overlapping, then you will see little enrichment and your experiment didn't work that well. If you're sample lines are switched – then you probably switched the sample names and we recommend doing the right thing and repeating the experiment and not simply switch sample names for the sake of a publication.



In this example, there is reasonable enrichment in our chip samples. And so we are confident we can see enrichment.

#### 2 – The Venn Diagram (well Venn Text)

In the folder “21\_peaks\_merging” folder, you will find the “venn.txt” file. This will show you a textual venn diagram of the overlap between the called peaks across all peak callers. To know your

experiment worked well, you should see a full list with combinations of all peak callers and relatively large numbers for the consensus peak sets (ie peaks called by multiple peak callers) – this is the ideal case. However, from our experience, there will almost always be 1 maybe 2 peak callers that don't like a dataset for some reason and so you may find a peak caller performed poorly but the others performed admirably. This is still a good and valid result. If you look at this file and only see small number of peaks and little overlap, and only 1 peak caller seems to have dominated peak

MACS2	SICER2	HOMER	Genrich	Total	Name
			X	103	Genrich
		X		2151	HOMER
		X	X	12	HOMER Genrich
X				14499	SICER2
X			X	328	SICER2 Genrich
X	X			10346	SICER2 HOMER
X	X	X	X	687	SICER2 HOMER Genrich
X				522	MACS2
X			X	606	MACS2 Genrich
X		X		78	MACS2 HOMER
X		X	X	44	MACS2 HOMER Genrich
X	X			1115	MACS2 SICER2
X	X		X	714	MACS2 SICER2 Genrich
X	X	X		12833	MACS2 SICER2 HOMER
X	X	X	X	28549	MACS2 SICER2 HOMER Genrich

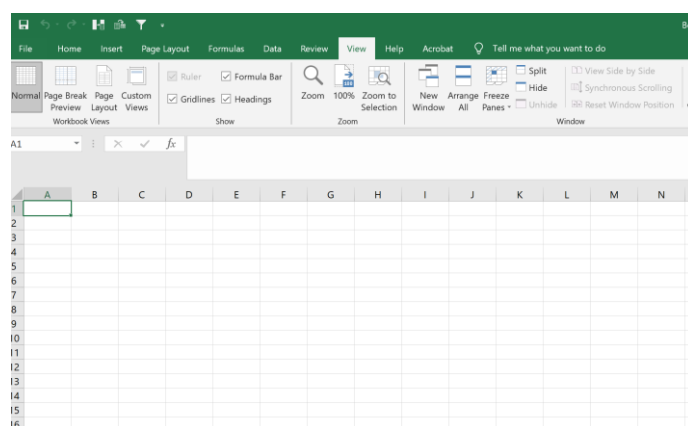


calling, then likely your experiment didn't work that great. Just because only 1 peak caller performed well though, doesn't mean the experiment is a write-off and a failure. It can still be valid and so doing some manual validations on the top fold-change differential peaks by chip-PCR might give you an indication whether there is salvageable data or not. Also if you have other confirmatory experimental evidence then even 1 peak caller getting results is fine. This is why we implemented multiple peak callers, because there are many instances where the signal:noise just creates a mess for most peak callers but generally 1 will be the super-hero of the day in such a situation.

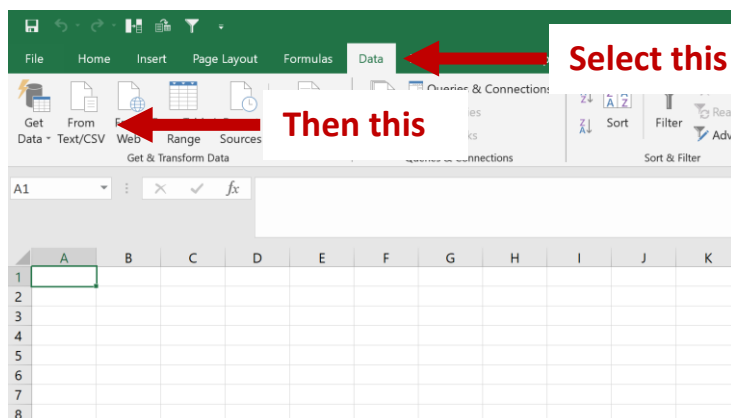
### 3 – What results files do I look at exactly?

Valid question. In the folder “22\_peak\_processing,” open the “xxxx\_all\_peaks\_calculated.tsv” file in excel and you're good to go. Now to open it there is a little step to do...

Open a new blank workbook in excel



In the ribbon at the top, go to “Data”, then select “From Text/CSV”



In the dialog box that opens up, find and open the peaks files “xxxx\_all\_peaks\_calculated.tsv.” Follow all the prompts and keep pressing “Next” / “Proceed” till the end and the file opens. Opening the peak file this way circumvents an issue that Excel constantly makes which is it will interpret some gene names such as OCT1 as a date, when its not. So by following the afore mentined steps, excel will not do this stupid conversion and instead, when you save the file as an xlsx, it will ensure that this issue doesn't happen (seen it in sooooo many publications its not funny – just import data this way please people?)

From this file, you can view all the results and data for you analysis. Refer to Interpreting ChIP-AP Output for the definition of what each column means.



#### 4 – How do I view my alignments and data?

People typically want to view their results on UCSC or other genome browsers. As we don't have a web-server to host such coverage files (and making accessible ucsc hub is a real pain and we don't want to implement that), the onus is on you to view them locally on your machine. All laptops, whether then can run ChIP-AP or not can run IGV [Downloads | Integrative Genomics Viewer \(broadinstitute.org\)](https://broadinstitute.org/Downloads/IntegrativeGenomicsViewer) and view the coverage and bam files. The coverage and bam files can be located in the "08\_results" folder.

Download IGV, install it (super easy) and then load the coverage and bam files needed. Make sure you load the right genome build however! That's critical. From section Main Pipeline Output, you can copy columns B,C,D straight into IGV and it will take you to the peak section.

#### 5 – In Short, whats relevant?

Easy answers

1 – check fingerprint plot and make sure it looks good

2 – check venn.txt file and make sure you get good spread of peaks

Together points 1 and 2 tell you your experiment worked!

3 – Your final peak file is in "22\_peak\_processing" open the "xxxx\_all\_peaks\_calculated.tsv" – This is the file you need to upload to GEO as your processed data file for your analysis and the only file you need to work with when looking through your data.

4 – Also as part of your submission to GEO or as a supplemental table in your manuscript, you MUST include the settings table named "default\_settings\_table.tsv" located in the root analysis directory. This provided with the raw fq files, which must be uploaded to GEO, will ensure complete reproducibility of the analysis performed.

5 – Manuscript details for M&M. A statement such as the following should suffice.

For processing our ChIP-Seq analysis, we utilized ChIP-AP (<https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbab537/6489109>). Raw fq files are uploaded to GEO with accession number GSE172355, and the custom settings table utilized for analysis can be found on GEO as a processed settings file and also in Table 1 in our manuscript. Full details of ChIP-AP and its function can be found in its corresponding manuscript (<https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbab537/6489109>).



## Manuals and Citations

If you use ChIP-AP in your analysis, please cite the us and all the following programs

Programs	References
ChIP-AP V5.2	Guide: <a href="https://github.com/JSuryatenggara/ChIP-AP/wiki/ChIP-AP-Guide">https://github.com/JSuryatenggara/ChIP-AP/wiki/ChIP-AP-Guide</a> GitHub: <a href="https://github.com/JSuryatenggara/ChIP-AP">https://github.com/JSuryatenggara/ChIP-AP</a> Citation: <a href="https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbab537/6489109">https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbab537/6489109</a>
Python3 Linux 3.7.x/3.8.x macOS 3.7.x	We have noted in our testing that there is a change in python 3.8 on macOS in how multi-threading is handled which breaks ChIP-AP. As such, for macOS installs you must ensure that python3.7.x is installed. If using our installation guides, the provided yml files will ensure all the correct dependencies and requirements are met automatically.
FastQC v0.11.9	Guide: <a href="https://www.bioinformatics.babraham.ac.uk/projects/fastqc/">https://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a> GitHub: <a href="https://github.com/s-andrews/FastQC">https://github.com/s-andrews/FastQC</a>
Clumpify v38.18 (BBmap)	Introduction: <a href="https://www.biostars.org/p/225338/">https://www.biostars.org/p/225338/</a> Guide: <a href="https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/clumpify-guide/">https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/clumpify-guide/</a> GitHub: <a href="https://github.com/BioInfoTools/BBMap/blob/master/sh/clumpify.sh">https://github.com/BioInfoTools/BBMap/blob/master/sh/clumpify.sh</a> Citation: <a href="https://www.osti.gov/biblio/1241166-bbmap-fast-accurate-splice-aware-aligner">https://www.osti.gov/biblio/1241166-bbmap-fast-accurate-splice-aware-aligner</a>
BBDuk v38.18 (BBmap)	Introduction: <a href="http://seqanswers.com/forums/showthread.php?t=42776">http://seqanswers.com/forums/showthread.php?t=42776</a> Guide: <a href="https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbduk-guide/">https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbduk-guide/</a> GitHub: <a href="https://github.com/BioInfoTools/BBMap/blob/master/sh/bbduk.sh">https://github.com/BioInfoTools/BBMap/blob/master/sh/bbduk.sh</a> Citation: <a href="https://www.osti.gov/biblio/1241166-bbmap-fast-accurate-splice-aware-aligner">https://www.osti.gov/biblio/1241166-bbmap-fast-accurate-splice-aware-aligner</a>
Trimmomatic v0.39	Guide: <a href="http://www.usadellab.org/cms/?page=trimmomatic">http://www.usadellab.org/cms/?page=trimmomatic</a> Downloadable manual page: <a href="http://www.usadellab.org/cms/uploads/supplementary/Trimmomatic/TrimmomaticManual_V0.32.pdf">http://www.usadellab.org/cms/uploads/supplementary/Trimmomatic/TrimmomaticManual_V0.32.pdf</a> GitHub: <a href="https://github.com/timflutre/trimmomatic">https://github.com/timflutre/trimmomatic</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4103590/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4103590/</a>
bwa v0.7.17	Guide: <a href="http://bio-bwa.sourceforge.net/bwa.shtml">http://bio-bwa.sourceforge.net/bwa.shtml</a> GitHub: <a href="https://github.com/lh3/bwa">https://github.com/lh3/bwa</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2705234/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2705234/</a>
samtools view v1.9 (samtools)	Guide: <a href="http://www.htslib.org/doc/samtools-view.html">http://www.htslib.org/doc/samtools-view.html</a> GitHub: <a href="https://github.com/samtools/samtools">https://github.com/samtools/samtools</a> Citation: <a href="https://pubmed.ncbi.nlm.nih.gov/19505943/">https://pubmed.ncbi.nlm.nih.gov/19505943/</a>
deeptools plotFingerprint v3.5.0 (deepTools)	Guide: <a href="https://deeptools.readthedocs.io/en/develop/content/tools/plotFingerprint.html">https://deeptools.readthedocs.io/en/develop/content/tools/plotFingerprint.html</a> Citation: <a href="https://academic.oup.com/nar/article/44/W1/W160/2499308?login=true">https://academic.oup.com/nar/article/44/W1/W160/2499308?login=true</a>
MACS2 v2.2.6	Guide: <a href="https://hbctraining.github.io/Intro-to-ChIPseq/lessons/05_peak_calling_mac.html">https://hbctraining.github.io/Intro-to-ChIPseq/lessons/05_peak_calling_mac.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2732366/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2732366/</a> GitHub: <a href="https://github.com/mac3-project/MACS/wiki">https://github.com/mac3-project/MACS/wiki</a>
GEM v2.7	Guide: <a href="https://groups.csail.mit.edu/cgs/gem/">https://groups.csail.mit.edu/cgs/gem/</a> GitHub: <a href="https://github.com/gifford-lab/GEM">https://github.com/gifford-lab/GEM</a> Citation: <a href="https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002638">https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002638</a>
SICER2 v1.0.2	Guide: <a href="https://zanglab.github.io/SICER2/">https://zanglab.github.io/SICER2/</a> GitHub: <a href="https://github.com/bioinf/SICER2">https://github.com/bioinf/SICER2</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2732366/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2732366/</a>



HOMER findPeaks v4.11 (HOMER)	Guide: <a href="http://homer.ucsd.edu/homer/ngs/peaks.html">http://homer.ucsd.edu/homer/ngs/peaks.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/</a>
Genrich v0.6	Guide: <a href="https://informatics.fas.harvard.edu/atac-seq-guidelines.html">https://informatics.fas.harvard.edu/atac-seq-guidelines.html</a> GitHub: <a href="https://github.com/jsh58/Genrich">https://github.com/jsh58/Genrich</a>
Homer mergePeaks v4.11 (HOMER)	Guide: <a href="http://homer.ucsd.edu/homer/ngs/mergePeaks.html">http://homer.ucsd.edu/homer/ngs/mergePeaks.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/</a>
HOMER annotatePeaks v4.11 (HOMER)	Guide: <a href="http://homer.ucsd.edu/homer/ngs/annotation.html">http://homer.ucsd.edu/homer/ngs/annotation.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/</a>
IDR v2.0.4.2	GitHub: <a href="https://github.com/nboley/idr">https://github.com/nboley/idr</a> Citation: <a href="https://projecteuclid.org/journals/annals-of-applied-statistics/volume-5/issue-3/Measuring-reproducibility-of-high-throughput-experiments/10.1214/11-AOAS466.full">https://projecteuclid.org/journals/annals-of-applied-statistics/volume-5/issue-3/Measuring-reproducibility-of-high-throughput-experiments/10.1214/11-AOAS466.full</a>
HOMER findMotifsGenome v4.11 (HOMER)	Guide: <a href="http://homer.ucsd.edu/homer/ngs/peakMotifs.html">http://homer.ucsd.edu/homer/ngs/peakMotifs.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2898526/</a>
MEME meme-chip V5.0.5 (MEME)	Guide: <a href="https://meme-suite.org/meme/doc/meme-chip.html">https://meme-suite.org/meme/doc/meme-chip.html</a> Citation: <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2703892/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2703892/</a>