

GENE-SWitCH – Protocols



GENE-SWitCH

The regulatory GENomE of SWine and CHicken: functional annotation during development

Protocol WP2 T2.2 Analysis of WGBS data

Authors: Jani de Vos(WU), Ole Madsen (WU)

Workpackage: WP2

Version: 1.0

Protocol associated with Deliverable(s):	D2.1, D2.2, D2.3
Submission date to FAANG:	27/01/2022
Means of verification:	N/A

Research and Innovation Action, SFS-30-2018-2019-2020 Agri-Aqua Labs Duration of the project: 01 July 2019 – 30 June 2023, 48 months



Table of contents

1	Sun	nmary	3
2	Pro	tocol description	3
	2.1	Running the pipeline for analysis of GENE-SWitCH WGBS data	3
	2.2	Merging WGBS data	4
	2.3	Processing merged bam files with pipeline v2	4
	2.4	Compiling bed output	4



1 Summary

GENE-SWitCH aims at identifying functional elements located in the genomes of pig and chicken working on seven different tissues at three different developmental stages. A multitude of genomic marks regulate gene expression, these marks work together in a complex manner to facilitate correct gene expression during a whole life cycle. These marks are: DNA methylation, protein modifications and non-coding RNA molecules. DNA methylation is the addition of a methyl group to cytosine which is typically located 5' from a guanine, which creates CpG sites throughout the genome in vertebrates. Bisulphite sequencing, is a method whereby unmethylated cytosine residuals are converted to uracils and methylated cytosines remain unchanged identifying methylation across the genome. Site specific DNA methylation changes are detected through this methodology.

The GSM-pipeline (GENE-SWitCH project methylation analysis) is based on the pre-existing methylseq nf-core pipeline (https://nf-co.re/methylseq/1.6.1) for analysis of methylation data to which the GENE-SWitCH team has added extensions (https://github.com/FAANG/GSM-pipeline).

Processes in this pipeline include:

- 1. **Quality control**: FastQC (https://github.com/s-andrews/FastQC)
- 2. **Trimming**: TrimGalore (http://www.bioinformatics.babraham.ac.uk/ projects/trim_galore/)
- 3. Alignment: Bismark(https://github.com/FelixKrueger/Bismark) or bwa-meth (<u>https://github.com/brentp/bwa-meth</u>)
- 4. **Methylation calling**: Bismark methyl caller, CGmaptools (https://github.com/guoweilong/cgmaptools), MethylDackel with bwa-meth aligner (https://github.com/dpryan79/methyldackel)
- Downstream analysis: CGmaptools (<u>https://github.com/guoweilong/cgmaptools</u>), MethylKit (<u>https://bioconductor.org/packages/release/bioc/html/methylKit.html</u>), viewBS (<u>https://github.com/xie186/ViewBS</u>)
- 6. Results output: MultiQC (<u>https://github.com/ewels/MultiQC</u>)

2 Protocol description

2.1 Running the pipeline for analysis of GENE-SWitCH WGBS data

Pipeline can be downloaded from Github (https://github.com/FAANG/GSM-pipeline/tree/version-1.0), and requires installation of newest version of nextflow, together with either singularity or docker. The pipeline (version 1.0) was launched in the following way for analysis of GENE-SWitCH single end raw reads:

nextflow run GSM-pipeline/ -profile singularity --input {raw_reads}.fq.gz --fasta {reference}.fa --aligner bismark

GENE-SWitCH – Milestone MS5



2.2 Merging WGBS data

Generated bam files from the pipeline was merged using samtools merge. Samples were merged as these were files from different lanes.

2.3 Processing merged bam files with pipeline v2

Merged bam files are provided as an input to pipeline v2 (https://github.com/FAANG/GSMpipeline/tree/version2.0), where deduplication is followed by methylation calling using bismark. Methylation report generated through bismark is converted to CGmap file and further downstream analysis is executed.

nextflow run GSM-pipeline/ -profile singularity --bam {file}.bam --fasta {reference}.fa --aligner none --skip_alignment

2.4 Compiling bed output

Methylation calling results can be found in results/cgmaptools/cgmap_methyl_call_CHR/ directory and the output {name}.CGmap file contains methylation levels per site. The file was compiled into a bed file, and filtered on a minimum of 10 reads per site, this parameter excludes non-informative reads.

Output BED files contains the following information:

- 1. Informative header
- 2. Columns:
- 1. CHROM=Chromosome
- 2. POS=Position of RRBS site
- 3. POS=Position of RRBS site +1
- 4. NUC=Nucleotide on reference genome
- 5. CONT=Context
- 6. DINUC=Dinucleotide context
- 7. METH=Methylation level
- 8. MC=Counts of reads that support methylated cytosine
- 9. NC=Counts of reads that support all cytosines