# target-methylseq-qc: a lightweight pipeline for collecting metrics from targeted sequence mapping files.

**Abhinav Sharma** [1], **Talya Conradie** [2,3], **David Martino** [2], **Stephen Stick** [2,4,5], **and Patricia Agudelo-Romero** [2,6,7]

**1** Division of Molecular Biology and Human Genetics, Faculty of Medicine and Health Sciences, Stellenbosch University, Cape Town. **2** Wal-yan Respiratory Research Centre, Telethon Kids Institute, WA, Australia **3** Medical, Molecular and Forensic Sciences, Murdoch University, WA, Australia **4** Department of Respiratory and Sleep Medicine, Perth Children's Hospital for Children, WA, Australia. **5** Centre for Cell Therapy and Regenerative Medicine, School of Medicine and Pharmacology, WA, Australia. **6** Australian Research Council Centre of Excellence in Plant Energy Biology, School of Molecular Sciences, The University of Western Australia, WA, Australia **7** European Virus Bioinformatics Center, TH, Germany.

## Summary

Next-generation targeted genome sequencing allows the analysis of regions of interest within a genome. While it is possible to incorporate targeted sequencing into whole-genome sequencing (WGS) bioinformatics pipelines, there remains a gap in accurately converting WGS metrics into precise target sequencing metrics and filtering the raw BED files into the targeted regions. Here, we introduce the target-methylseq-qc pipeline (https://doi.org/10.5281/zenodo.13147688), designed to (i) collect metrics from alignment files generated in targeted-methylation sequence analysis using the `picard_profiler` mode and (ii) filtering `bedGraph` for features overlapping with the reference BED file using the `bed_filter` mode, both of these modes are subworkflows written using the Nextflow (Di Tommaso et al., 2017) workflow language.

The target-methylseq-qc pipeline, when used in the `picard_profiler` mode accepts inputs in various alignment formats, including SAM, BAM and CRAM files (*HTS Format Specifications*, 2023). Additionally, to refine the metrics to the target regions the inclusion of a FASTA reference file and BED intervals file is required. Upon completion of the analysis, a MultiQC report (Philip Ewels et al., 2016) will be generated, encompassing the updated sequencing coverage data for the targeted regions with some extras. The `picard_profiler` mode of the pipeline integrates Picard metrics from GATK picard tools (McKenna et al., 2010; *Picard Toolkit*, 2019), using two specific metrics: (i) collecthsmetrics (*CollectHsMetrics (Picard)*, 2019), which relies upon the hybrid-selection technique to capture exon sequences for targeted sequencing experiments; and (ii) collectmultiplemetrics (*CollectMultipleMetrics (Picard)*, 2021), which captures closely related metrics such as alignment summary, insert size, and quality score.

On the other hand, `bed_filter` mode of the pipeline is designed to filter the bedGraph files outcome from nf-core/methylseq (Phil Ewels et al., 2024) using the reference bed panel, in this case the Twist Human Methylome panel (https://www.twistbioscience.com/resources/data-files/twist-human-methylome-panel-target-bed-file) and best practices *Twist Methylome* (2016b) using bedtools (Quinlan & Hall, 2010) filter command. Filtering raw BED files with the targeted regions is crucial because it ensures that the analysis focuses on specific genomic targets accurately and efficiently. This step minimizes the inclusion of off-target sequences and reduces the potential for including sequencing artifacts, which can be

44 introduced during capture-based targeted sequencing processes. Downstream analyses from
45 the filtered BED files will enable the calculation of CpG ratios and the testing for differentially
46 methylated cytosines (DMCs) or regions (DMRs).

47 Regardless of the usage mode of the pipeline, the final MultiQC report automatically collates
48 the relevant reports from FastQC (Andrews, 2010), Bedtool and Picard tools in an HTML
49 document, which could be shared with collaborators or added as supplementary material in
50 publications.

51 target-methylseq-qc is a portable pipeline compatible with multiple platforms, such as local
52 laptop or workstation machines, high-performance computing environments and cloud infra-
53 structure. Although target-methylseq-qc was originally created for calculating sequencing
54 coverage in target sequencing as a follow-up step to the nf-core/methylseq pipeline (Phil
55 Ewels et al., 2024), within the Airway Epithelium Respiratory Illnesses and Allergy (AERIAL)
56 paediatric cohort study (Kicic-Starcevich et al., 2023); its versatility allows for extending its
57 application to other sequencing panels from various next-generation methods.

## Statement of need

59 The target-methylseq-qc pipeline is designed to streamline the quality control process for target
60 methylation sequencing data. Researchers and bioinformaticians working with methylation
61 sequencing data often face challenges in ensuring data quality and consistency across different
62 samples and experiments. This pipeline addresses these challenges by providing a standardized
63 and automated workflow for quality control, leveraging the capabilities of the nf-core framework.
64 Key features of the target-methylseq-qc pipeline include (i) Standardized Input Format:
65 The pipeline expects a CSV samplesheet with specific fields tailored to different modes
66 (picard_profiler and bed_filter), ensuring consistency and ease of use (ii) Flexible Execution
67 Modes: Users can choose between different subworkflows (picard_profiler and bed_filter)
68 based on their specific needs, enabling tailored quality control processes (iii) Comprehensive
69 Parameter Control: Users can fine-tune the pipeline's behavior through a wide range of
70 parameters, covering execution modes, input/output options, reference genome options, and
71 infrastructural configuration. By automating and standardizing the quality control process, the
72 target-methylseq-qc pipeline helps researchers save time, reduce errors, and ensure high-quality
73 data for downstream analysis and clinically applicable insights.

## Design principles and capabilities

75 The target-methylseq-qc pipeline builds upon the standardised pipeline template maintained
76 by the nf-core community (P. A. Ewels et al., 2020) for Nextflow pipelines as well as makes
77 use of the nf-core/modules project to install modules for FastQC, MultiQC (Philip Ewels et
78 al., 2016), Bedtools, Picard as well as Samtools (Danecek et al., 2021) within the pipeline
79 Figure 1.

80 The use of the nf-core template facilitates keeping the design of the pipeline generic and
81 portable across different execution platforms, therefore the target-methylseq-qc pipeline can be
82 used on local machines, HPC orchestrators (e.g. SLURM, PBS), and cloud computing systems
83 such as AWS Batch, Azure Batch, Google Batch, in addition to the more generic Kubernetes
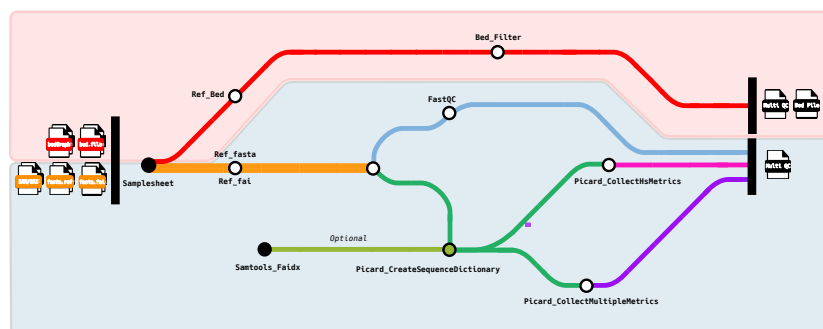84 distribution.

**Figure 1:** Subway map for various steps in the target-methylseq-qc pipeline.

In addition to the base workflow as mentioned in Figure 1, the pipeline also includes optional picard/createsequencedictionary (*CreateSequenceDictionary (Picard),* 2022) and Samtools modules to aid users in automatically generating the required genome dictionary (DICT) file, in case they have only the reference FASTA and BED files but intend to use the pipeline. Furthermore, depending on the quality check requirements of the users, we have enabled the metrics collection for 10x, 20x, 30x and 50x coverage.

## Tutorials and documentation

The steps needed to configure the pipeline inputs and configuration for the relevant infrastructure are available in the documentation within the GitHub repository as well as a dedicated documentation website (*Target-Methylseq-Qc Website*, 2024).

## Pre-requisites

To ensure proper operation of the target-methylseq-qc pipeline, three dependencies must be available in the execution environment: Java (LTS > 11), Nextflow (> 24.04), and a package manager such as conda (Gruning et al., 2018) or a container system such as docker or singularity (Veiga Leprevost et al., 2017).

Getting started with the pipeline setup is straightforward given that (i) Java (LTS > 11) (ii) Nextflow (> 24.04) and (iii) a package manager (e.g. conda) or a container system (e.g. docker or singularity) are available in the execution environment. The in-built test profile from the pipeline can then be used to execute the profile on the relevant infrastructure with some test dataset.

## Pipeline installation

target-methylseq-qc pipeline can be downloaded from the GitHub code repository using the git command line tool or directly through using the Nextflow command line tool using the following commands

```
# Git based download
$ git clone https://github.com/wal-yan/target-methylseq-qc
```

```
# Nextflow based download
$ nextflow pull https://github.com/wal-yan/target-methylseq-qc
```

## Test profiles

Two built-in test profiles are available in target-methylseq-qc pipeline for each mode of execution. These profiles can be used to run tests on the relevant infrastructure using the bundled test datasets (Agudelo-Romero, 2024), helping users to identify and resolve any infrastructural issue before the analysis stage.

```
# picard_profiler mode
$ nextflow run wal-yan/target-methylseq-qc \
  -profile docker,test_picard_profiler
```

```
# bed_filter mode
$ nextflow run wal-yan/target-methylseq-qc \
  -profile docker,test_bed_filter
```

## Input

Following the convention for standard input in the Nextflow pipelines, target-methylseq-qc expects a CSV samplesheet as an input with the following fields. An example of a samplesheet Table 1 for target-methylseq-qc in picard-profiler mode containing three columns, capturing the (i) name of the sample (ii) path to BAM file and (iii) path to the BAM index (BAI) file.

Table 1: Samplesheet structure for picard_profiler mode .

| sample | bam | bai |
|---|---|---|
| sample-01 | /path/to/sample-01.bam | /path/to/sample-01.bai |
| sample-02 | /path/to/sample-02.bam | /path/to/sample-02.bai |

Whereas the bed_filter mode requires a different set of columns in the input samplesheet CSV file, as shown in Table 2.

Table 2: Samplesheet structure for bed_filter mode .

| sample | bedGraph |
|---|---|
| sample-01 | /path/to/sample-01.bedGraph |
| sample-02 | /path/to/sample-02.bedGraph |

## Execution

The pipeline initialization step, as per the best practices of the nf-core template, checks the validity of the file paths specified to be either a POSIX-compliant file system or a cloud object storage path for files storaged in AWS S3, Azure Blob Storage or Google Cloud Storage buckets.

The behaviour of the pipeline can be controlled through the pipeline parameters which are divided into different groups such as (i) Execution Mode, (ii) Input/Output Options (iii) Reference Genome Options in addition to the generic parameters inherited from the nf-core

128 template such as (i) Max job request options (ii) Generic options and (iii) Institutional config
129 options. A complete list of the parameters specific to target-methylseq-qc pipeline is summarised
130 in Table 3

**Table 3:** Summary of pipeline-specific parameters for target-methylseq-qc pipeline .

| Parameter Name | Description |
| --- | --- |
| picard_profiler | Enable this boolean option to use the picard_profiler subworkflow |
| bed_filter | Enable this boolean option to use the bed_filter subworkflow |
| input | Path to comma-separated file containing information about the samples in the experiment. |
| outdir | The output directory where the results will be saved. |
| ref_fasta | Path to FASTA genome file. |
| ref_fai | Path to the FASTA index file. |
| ref_bed | Path to the BED file for the reference panel. |

## 131 Output

132 Upon completion, the two subworkflows generate different outputs which are presented together
133 in the MultiQC file. For picard_profile mode, a MultiQC file is produced, providing the relevant
134 results related to the coverage metrics Figure 2. For the bed_filter mode, a BED file is
135 generated with the methylation positions filtered based on the BED intervals file from the
136 targeted methylation profile Figure 3.

**Figure 2:** MultiQC report generated for target-methylseq-qc, in `picard-profiler` highlighting the refine metrics from targeted sequencing at 10X, 20X, 30X and 50X coverage.

**Figure 3:** Filtered bedGraph file generated using the `bed_filter` mode of target-methylseq-qc.

## References

Agudelo-Romero, P. (2024). *Wal-yan/target-methylseq-qc test-data* [Data set]. Zenodo. https://doi.org/10.5281/zenodo.13601364

Andrews, S. (2010). *Babraham bioinformatics - FastQC a quality control tool for high throughput sequence data*. https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

*CollectHsMetrics (picard). GATK*. (2019, November 25). https://gatk.broadinstitute.org/hc/en-us/articles/360036856051-CollectHsMetrics-Picard-

*CollectMultipleMetrics (picard). GATK*. (2021, February 22). https://gatk.broadinstitute.org/hc/en-us/articles/360057440491-CollectMultipleMetrics-Picard-

*CreateSequenceDictionary (picard). GATK*. (2022, November 12). https://gatk.broadinstitute.org/hc/en-us/articles/360036729911-CreateSequenceDictionary-Picard-

Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., Whitwham, A., Keane, T., McCarthy, S. A., Davies, R. M., & Li, H. (2021). Twelve years of SAMtools and BCFtools. *GigaScience*, *10*(2), giab008. https://doi.org/10.1093/gigascience/giab008

Di Tommaso, P., Chatzou, M., Floden, E. W., Barja, P. P., Palumbo, E., & Notredame, C. (2017). Nextflow enables reproducible computational workflows. *Nature Biotechnology*, *35*(4), 316–319. https://doi.org/10.1038/nbt.3820

Ewels, P. A., Peltzer, A., Fillinger, S., Patel, H., Alneberg, J., Wilm, A., Garcia, M. U., Di Tommaso, P., & Nahnsen, S. (2020). The nf-core framework for community-curated bioinformatics pipelines. *Nature Biotechnology*, *38*(3), 276–278. https://doi.org/10.1038/s41587-020-0439-x

Ewels, Phil, Hüther, P., Miller, E., Sateesh_Peri, Spix, N., bot, nf-core, Peltzer, A., F., S., Alneberg, J., Garcia, M. U., Krueger, F., Tommaso, P. D., Hörtenhuber, M., Ajith, V., Davenport, C., Patel, H., Salam, W., Cochetel, N., Alessia, … Céline, N. (2024). *Nf-core/methylseq: Huggy mollusc* (Version 2.6.0). Zenodo. https://doi.org/10.5281/zenodo.10463781

Ewels, Philip, Magnusson, M., Lundin, S., & Käller, M. (2016). MultiQC: Summarize analysis

167 results for multiple tools and samples in a single report. *Bioinformatics*, *32*(19), 3047–3048.
168 https://doi.org/10.1093/bioinformatics/btw354

169 Gruning, B., Dale, R., Sjodin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., Valieris,
170 R., Koster, J., & Bioconda, T. (2018). Bioconda: Sustainable and comprehensive software
171 distribution for the life sciences [Journal Article]. *Nat Methods*, *15*(7), 475–476. https:
172 //doi.org/10.1038/s41592-018-0046-7

173 *HTS format specifications*. (2023). https://samtools.github.io/hts-specs/

174 Kicic-Starcevich, E., Hancock, D. G., Iosifidis, T., Agudelo-Romero, P., Caparros-Martin, J.
175 A., Silva, D., Turkovic, L., Souef, P. N. L., Bosco, A., Martino, D. J., Kicic, A., Prescott,
176 S. L., & Stick, S. M. (2023). *Airway epithelium respiratory illnesses and allergy (AERIAL)*
177 *birth cohort: Study protocol*. medRxiv. https://doi.org/10.1101/2023.04.29.23289314

178 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella,
179 K., Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The genome analysis
180 toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data.
181 *Genome Research*, *20*(9), 1297–1303. https://doi.org/10.1101/gr.107524.110

182 *Picard toolkit*. (2019). https://broadinstitute.github.io/picard/

183 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing ge-
184 nomic features. *Bioinformatics*, *26*(6), 841–842. https://doi.org/10.1093/bioinformatics/
185 btq033

186 *Target-methylseq-qc website*. (2024). https://wal-yan.github.io/target-methylseq-qc/usage.
187 html

188 *Twist methylome*. (2016a). https://www.twistbioscience.com/products/ngs/fixed-panels/
189 human-methylome-panel

190 *Twist methylome*. (2016b). https://www.twistbioscience.com/resources/technical-note/
191 analyzing-twist-targeted-methylation-sequencing-data-using-twist-human

192 Veiga Leprevost, F. da, Gruning, B. A., Alves Aflitos, S., Rost, H. L., Uszkoreit, J., Barsnes,
193 H., Vaudel, M., Moreno, P., Gatto, L., Weber, J., Bai, M., Jimenez, R. C., Sachsenberg,
194 T., Pfeuffer, J., Vera Alvarez, R., Griss, J., Nesvizhskii, A. I., & Perez-Riverol, Y. (2017).
195 BioContainers: An open-source and community-driven framework for software standard-
196 ization [Journal Article]. *Bioinformatics*, *33*(16), 2580–2582. https://doi.org/10.1093/
197 bioinformatics/btx192