

1 cfDNA UniFlow: A unified preprocessing pipeline for cell-free DNA
2 data from liquid biopsies

3

4 **Authors**

5 Sebastian Röner¹, Lea Burkard^{1,4}, Michael R. Speicher^{2‡}, Martin Kircher^{1,3}

6 ¹ Berlin Institute of Health (BIH) at Charité – Universitätsmedizin Berlin, Berlin, Germany

7 ² Institute of Human Genetics, Diagnostic and Research Center for Molecular BioMedicine, Medical
8 University of Graz, Graz, Austria

9 ³ Institute of Human Genetics, University Medical Center Schleswig-Holstein, University of Lübeck,
10 23562 Lübeck, Germany

11 ⁴ University of Potsdam, Institute for Biochemistry and Biology, 14469 Potsdam, Germany

12 ‡ Deceased on Sep 24, 2023

13

14 *Contact:* martin.kircher@bih-charite.de

15

16 Abstract:

17 **Background:**

18 Cell-free DNA (cfDNA), a broadly applicable biomarker commonly sourced from urine or blood, is
19 extensively used for research and diagnostic applications. In various settings, genetic and epigenetic
20 information is derived from cfDNA. However, a unified framework for its processing is lacking, limiting
21 the universal application of innovative analysis strategies and the joining of data sets.

22 **Findings:**

23 Here, we describe cfDNA UniFlow, a unified, standardized, and ready-to-use workflow for processing
24 cfDNA samples. The workflow is written in Snakemake and can be scaled from stand-alone computers
25 to cluster environments. It includes methods for processing raw genome sequencing data as well as
26 specialized approaches for correcting sequencing errors, filtering, and quality control. Sophisticated
27 methods for detecting copy number alterations and estimating and correcting GC-related biases are
28 readily incorporated. Furthermore, it includes methods for extracting, normalizing and visualizing
29 coverage signals around user defined regions in case-control settings. Ultimately, all results and
30 metrics are aggregated in a unified report, enabling easy access to a wide variety of information for
31 further research and downstream analysis.

32 **Conclusions:**

33 We provide an automated pipeline for processing cell-free DNA sampled from liquid biopsies, including
34 a wide variety of additional functionalities like bias correction and signal extraction. With our focus on
35 scalability and extensibility, we provide a foundation for future cfDNA research and faster clinical
36 applications.

37

38 **Keywords:**

39 Cell-free DNA, liquid biopsies, sequence analysis, cancer detection, workflow

40

41 **Issue Section:**

42 Technical Note

43

44 **Availability and implementation:** Source code and extensive documentation is available on our
45 GitHub repository (<https://github.com/kircherlab/cfDNA-UniFlow>).

46

47

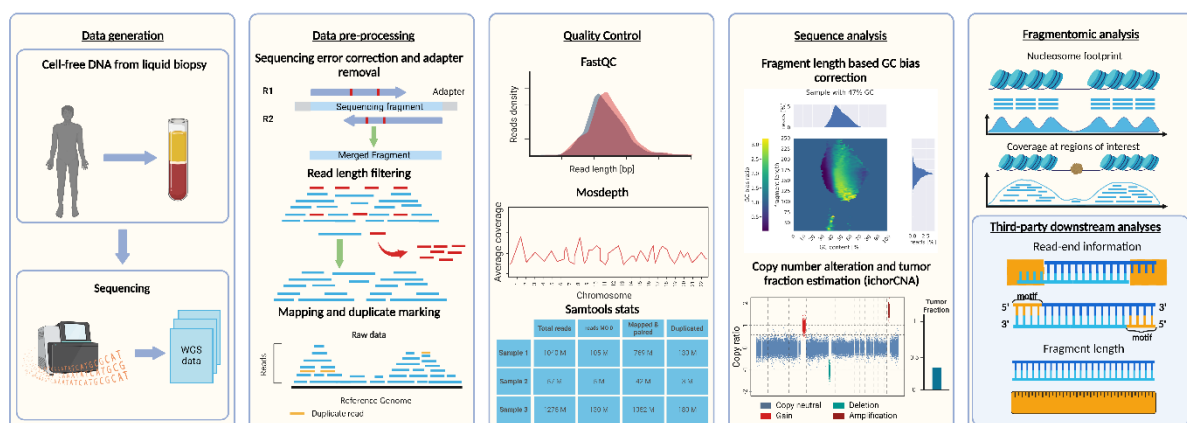
48 Introduction/Background

49 Cell-free DNA (cfDNA) is found in many bodily fluids like blood plasma and urine [1]. It is believed to
 50 be primarily derived from natural degradation processes during cell turnover [2]. However, the
 51 proportion of cell-types and tissues contributing to cfDNA changes in the context of certain
 52 physiological conditions or disease processes [3,4]. Thus, signals in cfDNA might serve as relevant
 53 biomarkers in health and disease. Collecting cfDNA in so-called liquid biopsies (Fig. 1) is considered
 54 non-invasive and led to an increased research interest in the biomedical field for using cfDNA in
 55 allograft (i.e., donor organ) rejection, prenatal testing and diagnostics, as well as disease detection and
 56 health monitoring [5] (especially for cancer).

57 Over the last years, many approaches have been developed to extract information from cfDNA samples
 58 for various DNA applications. Methods range from identifying allelic differences at known disease markers,
 59 detection and tracking of mutations [6] and copy number alterations (CNAs) in tumor cells [6], and DNA
 60 fragmentation differences [3,8,9] to measuring methylation state [10–12]. While these methods
 61 exploit different signals, all rely on the precise quantification of read distributions, and slight changes
 62 in read recovery affect their results (Fig. 1).

63 Therefore, consistent data quality is the primary requirement for developing these new diagnostic
 64 methods (Fig. 1). Even though sample handling is constantly streamlined, individual differences of
 65 sample donors, and logistic factors like time of sample collection, duration, conditions of storage, and
 66 further preanalytical handling are challenging to fully control in a clinical context, but have been shown
 67 to affect the quality of cfDNA samples [13–16]. Additionally, detecting signals of interest (e.g., from
 68 circulating tumor DNA, ctDNA) in a background mainly derived from hematopoietic cells [16] is not
 69 trivial, emphasizing the need for optimal data quality.

70 One way to mitigate some preanalytical effects and technical biases introduced during sequencing of
 71 cfDNA samples is to include specialized correction and sampling steps during computational processing
 72 of the data (Fig. 1). Even though the need has been identified previously in the field of cfDNA,
 73 community standards are still lacking for preprocessing genome sequencing data from cfDNA [7,8,18–
 74 23].



75

76 *Figure 1: Overview of cfDNA analysis. The leftmost panel depicts data generation by liquid biopsy sampling followed by library
 77 preparation and sequencing. The second panel shows the entry point of cfDNA Uniflow. It displays the core functionality of
 78 merging reads/removing adapters, length filtering, mapping to a reference genome and duplicate marking. Sample quality
 79 control is shown in the third panel and for example performed using FastQC, Mosdepth and SAMtools stats. The fourth panel
 80 shows optional steps of GC bias correction and estimation of copy number alterations and tumor proportion. Finally, results
 81 are aggregated, for example in a report and used for downstream analyses (fifth panel). Figure created with Biorender.*

82 One reason for the lack of dedicated cfDNA pipelines, might be that many publications in the field are
83 focused on the downstream analysis like the classification of disease samples, relying on unpublished
84 in-house pipelines for data processing. Further, important correction steps are often tailored towards
85 specific features and tightly integrated in downstream analysis pipelines, making it difficult to
86 generalize and transfer them to new projects [7,8]. Nevertheless, there have been some approaches
87 trying to address the need for community standards. A notable one is the FinaleDB project, which
88 aggregates cfDNA samples from multiple sources, processes them in a uniform manner and provides
89 fragment coordinates via a web portal [23]. To protect the privacy of patients, the data is anonymized
90 during processing, removing all sequence information and making it unsuitable for analyses not
91 focused on fragmentation patterns. Additionally, this pipeline does not address issues of batch and
92 bias correction, which might be most relevant in such aggregation efforts. The project getting closest
93 to setting a community standard for processing samples not just for fragmentomics applications, is
94 called cfDNApipe. It combines many useful tools for basic processing of normal and bisulfite converted
95 DNA sequences. The utility functions range from generation of summary statistics, GC-bias correction
96 tailored towards CNV detection and extraction of a limited number of features [25]. However, the
97 software seems to be designed for single computer use, lacking many of the features provided by a
98 full-fledged workflow management system, making it hard to scale analysis in different environments,
99 like compute clusters. Moreover, the design does not allow for easy integration of new functionalities,
100 creating the need for either an additional workflow management system or extensive modification of
101 the original code.

102 Technical biases and missing community standards cause several drawbacks for the field. First, users
103 rely on standard processing pipelines from other fields, which might not be suitable for specific
104 analyses. They might also feel the need to develop their own pipelines by selecting appropriate tools
105 and tuning parameters optimized on the available set of samples. Second, it adds additional overhead
106 when comparing across multiple studies. Here, researchers are frequently required to work with the
107 original processing of each site, potentially introducing technical biases in the analysis. Alternatively,
108 reprocessing data from multiple sites can reduce technical biases between studies but creates an
109 additional computational and organizational burden (incl. access to raw and protected genetic data).
110 Third, it can be hard to keep track of all sample-level information when building analysis pipelines using
111 many samples, mainly when information gets scattered across many samples and files.

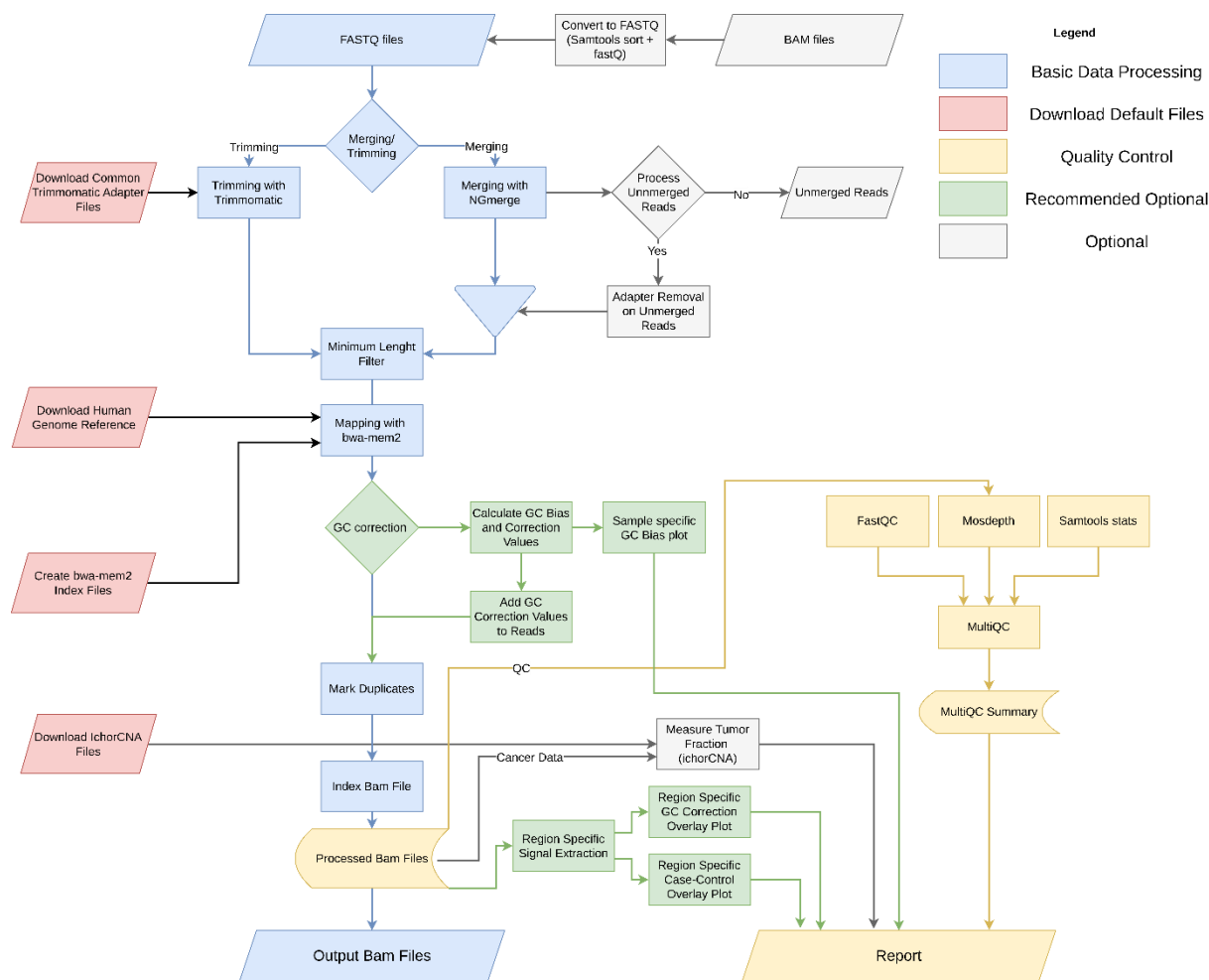
112 To jointly address several of these problems, we developed an easy-to-use unified preprocessing
113 workflow for cell-free DNA written in Snakemake. It combines a curated list of tools for processing
114 genomic cfDNA samples, custom tools for reducing technical biases, and tools for estimating additional
115 characteristics like copy number states. Our pipeline is implemented with high configurability,
116 scalability from single computers to high-performance compute clusters, and a sophisticated reporting
117 system.

118

119 Overview and implementation

120 Implementation

121 We implemented the cfDNA UniFlow workflow in the popular workflow management system
 122 Snakemake [25]. This makes it easy to scale the workflow in different computing environments and
 123 allows for parallel processing of multiple samples. Further, most of the rules are implemented to
 124 enable multiprocessing and efficiently utilize multiple cores for each task. Conveniently, default
 125 resources like genome references or standard adapter files can be downloaded, if not configured to
 126 point to already available resources. A detailed overview of the workflow is available in Figure S1.
 127 Briefly, cfDNA UniFlow covers three parts between data generation and downstream analysis: data pre-
 128 preprocessing, quality control and utility functions (Figure 1).



129
 130 *Figure 2: Overview of unified cfDNA preprocessing workflow. Functionalities are color coded by task. Blue boxes contain the*
 131 *core functionality of cfDNA Uniflow. Red boxes represent rules for the automatic download of public resources. Yellow boxes*
 132 *summarize the Quality Control and reporting steps. Finally, grey and green boxes are optional steps, with green boxes being*
 133 *highly recommended.*

134

135 Preprocessing

136 The core preprocessing steps (Fig 2., components depicted in blue) expect FASTQ files as input.
 137 Alternatively, existing alignments (BAM files) can be provided for reprocessing. In the latter case, the
 138 workflow automatically converts these to FASTQ files using SAMtools [26]. Afterwards, reads can be
 139 merged with NGmerge [27], which also removes sequencing library adapters and corrects sequencing
 140 errors and ambiguous bases based on the read-overlap consensus. Reads that were not merged, can

141 be postprocessed using NGmerge adapter removal mode and can be included in the mapping process.
142 Alternatively, the merging step can be skipped and reads will only be trimmed using Trimmomatic.
143 Prior to mapping with bwa-mem2 [29], reads are further filtered based on their length, excluding reads
144 that are shorter than a configurable threshold. Finally, duplicate reads are marked (SAMtools markdup)
145 before the BAM files of the samples are passed to the next step.

146 Quality Control

147 In the Quality control (QC) step (Fig. 2, components depicted in yellow), general post-alignment
148 statistics and graphs are calculated for each sample with SAMtools stats [26] and FastQC [30].
149 Additional information on sample-wide median coverage and coverage at different genomic regions is
150 calculated via Mosdepth [31]. The QC results are aggregated in an HTML report via MultiQC [32] and
151 an example is shown Figure S2.

152 Signal Extraction

153 In the last step, additional utility modules (Fig. 2, components depicted in green) can be configured
154 and executed. This includes our in-house GC bias estimation and correction methods
155 (https://github.com/kircherlab/cfDNA_GCcorrection), an extension of the method described by
156 Benjamini & Speed [17]. As fragmentation in cfDNA is driven by natural degradation processes, libraries
157 constructed from liquid biopsies tend to have fragments of a wide range of lengths and do not follow
158 the original assumption that length is well-approximated by the mean fragment length. Therefore, we
159 estimate the expected fragment distribution by sampling regions along the reference genome,
160 counting all possible fragments for a specified range of fragment lengths and sorting them in bins of
161 their GC content. Afterwards, we measure the sample specific fragment distribution in the same
162 regions, scale them and compare them to the theoretical distribution. Based on the ratio of observed
163 and expected, we calculated correction values for each fragment length and GC content. The resulting
164 weights are attached to the reads as tags, which can be used for a wide variety of downstream signal
165 extraction methods, while preserving the original read coverage and fragmentation patterns. We
166 provide specialized signal extraction routines to extract coverage derived signals using read weights.
167 Further, we included the widely used tool ichorCNA [6], to identify copy number alterations and
168 estimate tumor fraction. An example of the output is available in Figure S3.

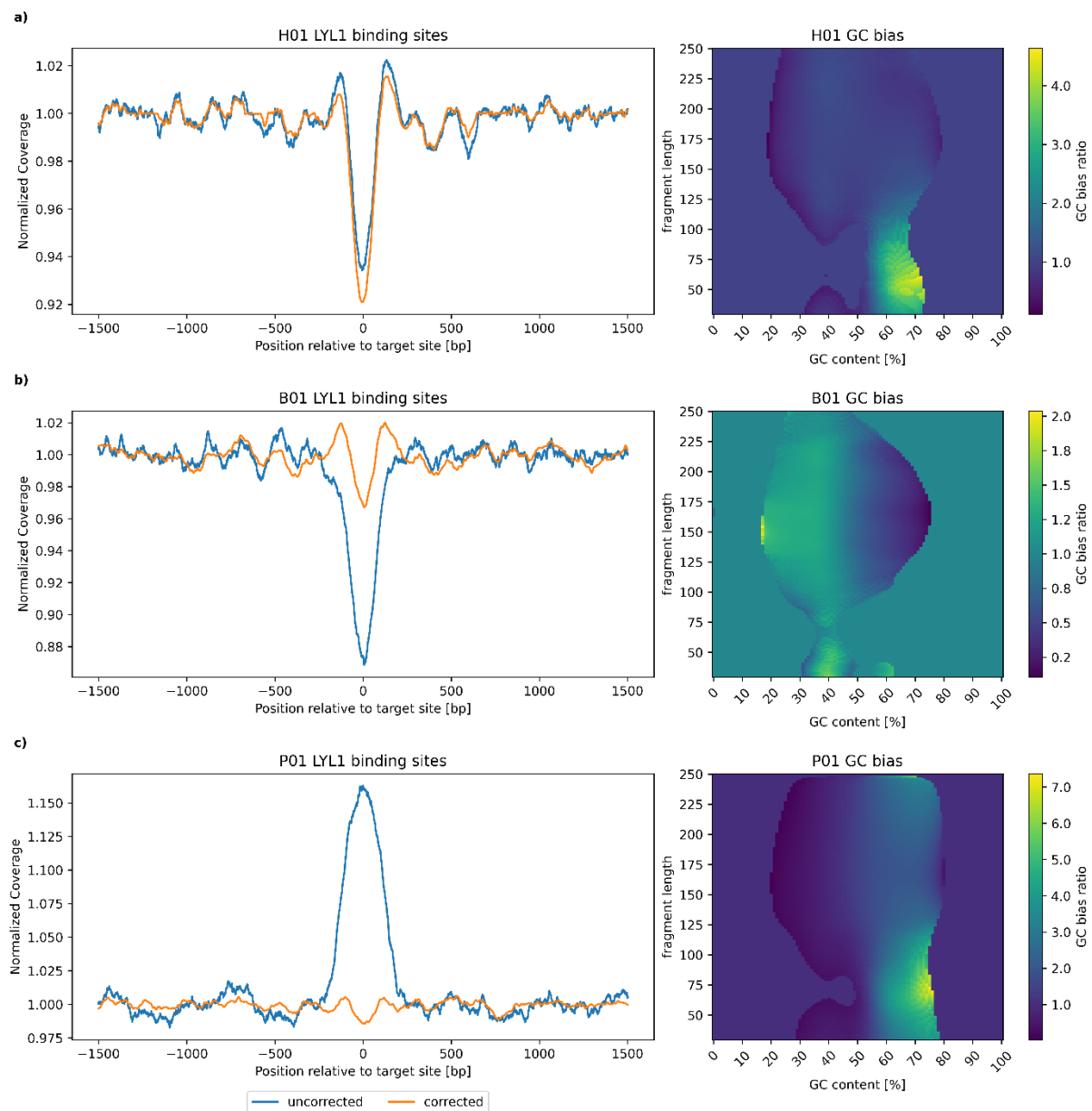
169 Reporting

170 Finally, all information provided by the previous steps is aggregated in a comprehensive HTML report.
171 This includes summary statistics on workflow execution provided by Snakemake, and plots and
172 summary statistics produced in the quality control steps. Additional information from optional steps
173 includes a general estimation of sample-specific GC bias parameters (Figure S4), the effects of GC bias
174 correction in user defined regions (Figure S5) and plots on copy number alterations created by
175 ichorCNA. Finally, case-control plots are generated and included, if more than one class of samples is
176 provided (Figure S6).

177 Results

178 To test and showcase cfDNA Uniflow, we use three exemplary cfDNA samples (healthy H01, breast
179 cancer B01, prostate cancer P01) with different conditions and average GC contents from the European
180 Genome-Phenome Archive Study EGAS00001006963. Each sample was converted to FASTQ files and
181 processed in our pipeline with standard parameters for human reference build GRCh38/hg38. As user-
182 defined regions of interest, we selected 10,000 binding sites of LYL1, a transcription factor (TF)
183 associated with hematopoietic cells [33], and GRHL2, an important pioneer TF for epithelial cells [34–
184 36] playing a role in a wide variety of cancer types [37–41]. Both TFs are especially suited due to their
185 association with expected tissue contributions in our samples and because they have high GC content
186 binding sites.

187 This can be seen in Figure 3, which shows coverage overlays centered on LYL1 binding sites and
188 illustrates the global and regional effects of GC biases in the respective samples. The healthy sample
189 H01, with an average GC content of 45%, shows a balanced global GC profile (Fig. 3a) and, accordingly,
190 the GC bias correction shows almost no effects on the composite signal. We see the strongest drop of
191 coverage at the TF binding site, expected for a sample of mainly hematopoietic origin where many
192 LYL1 binding sites are expected to be accessible to the TF. B01, a breast cancer sample with an average
193 GC content of 38%, shows an overrepresentation of fragments with GC content lower than the genome
194 average and an underrepresentation of fragments with higher GC content (Fig. 3b). This leads to a
195 distortion of the composite coverage signal around the LYL1 binding sites. Without GC correction, the
196 drop in coverage would be overestimated. After correction, coverage at the site is closer to the
197 coverage of the surrounding regions, consistent with an expected signal dilution compared to the
198 healthy sample (Fig. 3a) due to a higher contribution of non-hematopoietic cell-types in this cancer
199 sample (Fig. 3b). The same should be true for sample P01 (Fig. 3c), a prostate cancer sample with an
200 average GC content of 45%. However, the global GC bias profile (right panels) show the inverse trend
201 to sample B01, with a shift of fragment distribution towards a higher GC content. Unsurprisingly, the
202 signal around the binding sites is distorted towards higher coverage prior to the GC correction (i.e.,
203 suggesting that the TF binding sites are not accessible). After GC correction, the signal looks similar to
204 the one shown for B01, less open than the healthy sample and consistent with an increased
205 contribution of non-hematopoietic cell-types. Global effects of GC bias correction on fragment
206 distribution and a comparison to two other fragment-based are provided in the Supplement (Figure
207 S7).

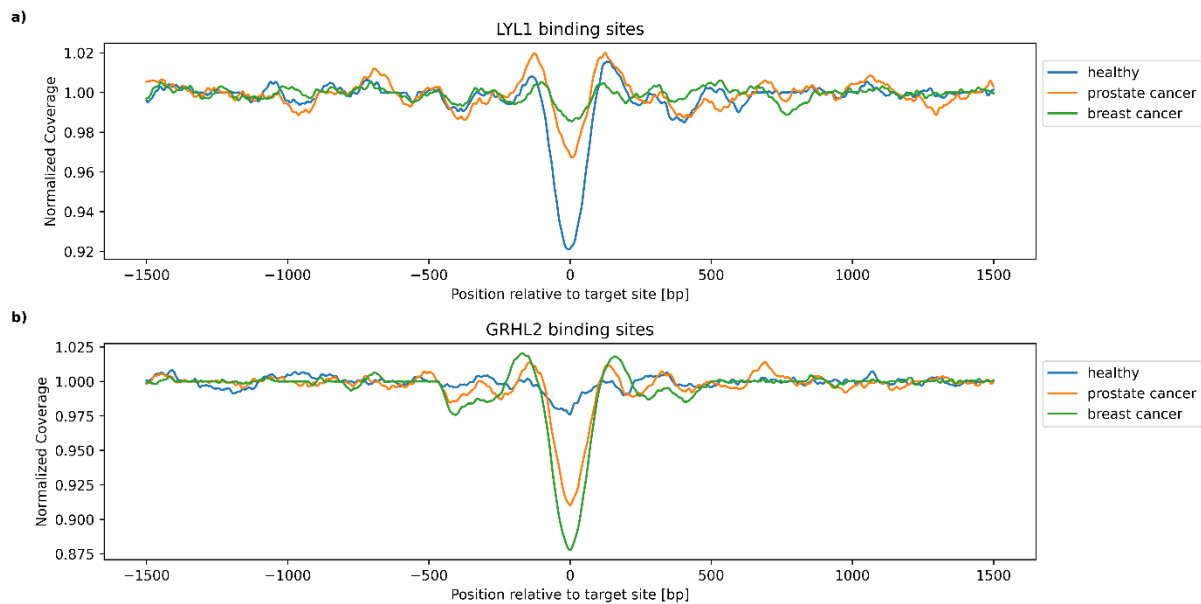


208

209 *Figure 3: Effects of GC bias on regional and global scale. Composite coverage signals of 10,000 LYL1 transcription factor*
 210 *binding sites (left) and global bias profiles (right) for three cfDNA samples are shown. a) Signals and profile for a healthy*
 211 *sample (H01) with an average GC content of 41%. The GC profile (right) is relatively balanced between observed and expected*
 212 *fragments. Respectively, the GC bias corrections have only minor effects on the composite coverage signal (left). b) GC bias*
 213 *effects for a breast cancer sample with an average GC content of 38%. The global GC profile shows an overrepresentation of*
 214 *lower GC content fragments (brighter color) and an underrepresentation of higher GC content fragments (darker color). This*
 215 *results in an underestimation of coverage (overestimation of accessibility) at the LYL1 binding sites. After GC correction the*
 216 *signal is closer to the surrounding coverage, consistent with lower relative contributions of hematopoietic cells and fewer open*
 217 *sites. In contrast, c) shows the GC bias effects of a prostate cancer cell with an average GC content of 45%. The global GC*
 218 *profile is skewed towards a higher GC content, leading to an overestimation of coverage around the LYL1 binding sites. After*
 219 *GC correction, the signal is closer to the surrounding coverage, indicating lower contributions of hematopoietic cells with*
 220 *accessible LYL1 sites.*

221 In addition to the GC bias plots for individual samples, we provide case-control plots for comparing
 222 sample classes with a control in the same plot. In our example, the healthy sample H01 would be the
 223 control and we are comparing samples on the LYL1 and GRHL2 sites. As noted, the expected signal
 224 around LYL1 binding sites is a drop in coverage for samples mainly derived from hematopoietic cells.
 225 When the contribution of non-hematopoietic cell-types, in which LYL1 is not expressed, increases, we
 226 expect to see a relative increase in coverage around the binding sites. Accordingly, the signals shown
 227 for our three test samples (Fig. 4a) are in line with that expectation. For GRHL2, we expect the opposite

228 signal. As healthy samples should not include many contributions from tissues with high GRHL2
229 activity, the expected coverage signal should be similar to the surrounding regions. In contrast, samples
230 with high contributions of cancer-derived DNA should show a drop in coverage, indicative of higher
231 accessibility of the TF binding sites (Fig. 4b).



232
233 *Figure 4: Case-control plots around GRHL2 and LYL1 binding sites. a) GC corrected composite coverage signals around 10,000*
234 *centered LYL1 transcription factor binding sites. The healthy sample (H01) shows lower relative coverage (i.e., higher*
235 *accessibility) at the center of the binding site overlay. This is consistent with higher LYL1 activity in hematopoietic cells. In*
236 *contrast, both cancer samples show higher relative coverage in the central region, in line with a higher proportion of non-*
237 *hematopoietic cells contributing to the signal. b) Composite coverage signals around 10,000 GRHL2 binding sites after GC*
238 *correction. Both cancer samples show a lower central coverage compared to surrounding regions (i.e., higher accessibility),*
239 *indicating higher activity than in the healthy sample. This is consistent with GRHL2 expression being associated with different*
240 *cancers.*

241 As pointed out before, exemplary figures of the other report sections, like QC or ichorCNA, can be
242 found in the supplement (Figures S2-S6). The full example report can be found in our GitHub
243 repository.

244

245 Conclusion

246 Here we propose cfDNA UniFlow, a unified preprocessing pipeline specifically tailored for cfDNA
247 samples. It is an easy-to-use, scalable, and configurable workflow, aiming to set a community standard
248 for enabling accessible and easily sharable future research in the field. In designing our workflow, we
249 aimed at providing a tool that can be used without much computer science background, but with the
250 option to be easily extended by experienced users with their own custom modules, allowing its
251 extension from a standard processing workflow to a full-featured analysis pipeline.

252 Availability and Requirements

253 Project name: [cfDNA-UniFlow](#)
254 Project home page: <https://github.com/kircherlab/cfDNA-UniFlow>
255 Operating system(s): Linux (64-bit)
256 Programming language: Python
257 Other requirements: Mamba or Conda
258 License: MIT

259

260 Additional Files

261 Supplementary Note: A general overview of the workflow including Supplementary Figures and a
262 quick start guide.

- 263 • Supplementary Figure 1: A detailed overview of the workflow components.
- 264 • Supplementary Figure 2: Example section of a QC report.
- 265 • Supplementary Figure 3: Example report for ichorCNA plot of copy number alterations.
- 266 • Supplementary Figure 4: Example of a sample's global GC bias estimate.
- 267 • Supplementary Figure 5: Example report of regional effects of GC bias correction.
- 268 • Supplementary Figure 6: Example of a case-control plot.
- 269 • Supplementary Figure 7: Comparison of three fragment based GC correction methods.

270

271 Abbreviations

272 BAM: binary alignment map; cfDNA: cell-free DNA; CNA: copy number alteration; ctDNA: circulating
273 tumor DNA; TF: transcription factor; QC: quality control

274

275 Data Availability

276 The data used for generating plots included in this article are available in the European Genome-
277 Phenome Archive at <https://ega-archive.org> with accession [EGAD00001010100](https://ega-archive.org).

278

279 Acknowledgements

280 We thank current and previous members of the Kircher and Speicher laboratories for helpful
281 discussions and suggestions. Computation has been performed on the HPC for Research cluster of the
282 Berlin Institute of Health at Charité – Universitätsmedizin Berlin.

283

284 Author contributions

285 Conceptualization: S.R., M.K. and M.R.S.; Data curation: S.R.; Formal analysis: S.R.; Funding acquisition:
286 M.K.; Methodology: S.R.; Project administration: M.K.; Resources: M.K. and M.R.S.; Software: S.R.;
287 Supervision: M.K.; Validation: S.R. and L.B.; Visualization: S.R.; Writing – original draft: S.R.; Writing -
288 review & editing: S.R., L.B., M.R.S. and M.K.

289

290 Conflict of interest

291 None declared.

292

293 References

- 294 1. Chan AKC, Chiu RWK, Lo YMD, Clinical Sciences Reviews Committee of the Association of Clinical
295 Biochemists. Cell-free nucleic acids in plasma, serum and urine: a new tool in molecular diagnosis.
296 *Ann Clin Biochem.* 2003; doi: 10.1258/000456303763046030.
- 297 2. Lo YM, Zhang J, Leung TN, Lau TK, Chang AM, Hjelm NM. Rapid clearance of fetal DNA from
298 maternal plasma. *Am J Hum Genet.* 1999; doi: 10.1086/302205.
- 299 3. Snyder MW, Kircher M, Hill AJ, Daza RM, Shendure J. Cell-free DNA Comprises an In Vivo
300 Nucleosome Footprint that Informs Its Tissues-Of-Origin. *Cell.* 2016; doi: 10.1016/j.cell.2015.11.050.
- 301 4. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al.. Inference of transcription factor binding
302 from cell-free DNA enables tumor subtype prediction and early detection. *Nat Commun.* Nature
303 Publishing Group; 2019; doi: 10.1038/s41467-019-12714-4.
- 304 5. Ding SC, Lo YMD. Cell-Free DNA Fragmentomics in Liquid Biopsy. *Diagn Basel Switz.* 2022; doi:
305 10.3390/diagnostics12040978.
- 306 6. Tunc I, Agbor-Enoh S, Valantine H, Thein SL, Pirooznia M. Cfcloud: A Cloud-Based Workflow for
307 Cell-Free DNA Data Analysis. *Blood.* 2020; doi: 10.1182/blood-2020-138785.
- 308 7. Adalsteinsson VA, Ha G, Freeman SS, Choudhury AD, Stover DG, Parsons HA, et al.. Scalable whole -
309 exome sequencing of cell-free DNA reveals high concordance with metastatic tumors. *Nat Commun.*
310 2017; doi: 10.1038/s41467-017-00965-y.
- 311 8. Peneder P, Stütz AM, Surdez D, Krumbholz M, Semper S, Chicard M, et al.. Multimodal analysis of
312 cell-free DNA whole-genome sequencing for pediatric cancers with low mutational burden. *Nat*
313 *Commun.* 2021; doi: 10.1038/s41467-021-23445-w.
- 314 9. Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al.. Genome-wide cell-free DNA
315 fragmentation in patients with cancer. *Nature.* 2019; doi: 10.1038/s41586-019-1272-6.
- 316 10. Erger F, Nörling D, Borchert D, Leenen E, Habbig S, Wiesener MS, et al.. cfNOME — A single assay
317 for comprehensive epigenetic analyses of cell-free DNA. *Genome Med.* 2020; doi: 10.1186/s13073-
318 020-00750-5.
- 319 11. Shen SY, Singhanian R, Fehringer G, Chakravarthy A, Roehrl MHA, Chadwick D, et al.. Sensitive
320 tumour detection and classification using plasma cell-free DNA methylomes. *Nature.* Nature
321 Publishing Group; 2018; doi: 10.1038/s41586-018-0703-0.
- 322 12. Chen S, Petricca J, Ye W, Guan J, Zeng Y, Cheng N, et al.. The cell-free DNA methylome captures
323 distinctions between localized and metastatic prostate tumors. *Nat Commun.* Nature Publishing
324 Group; 2022; doi: 10.1038/s41467-022-34012-2.
- 325 13. Jung M, Klotzek S, Lewandowski M, Fleischhacker M, Jung K. Changes in concentration of DNA in
326 serum and plasma during storage of blood samples. *Clin Chem.* 2003; doi: 10.1373/49.6.1028.
- 327 14. Lampignano R, Neumann MHD, Weber S, Klotten V, Herdean A, Voss T, et al.. Multicenter
328 Evaluation of Circulating Cell-Free DNA Extraction and Downstream Analyses for the Development of
329 Standardized (Pre)analytical Work Flows. *Clin Chem.* 2020; doi: 10.1373/clinchem.2019.306837.
- 330 15. Parpart-Li S, Bartlett B, Popoli M, Adleff V, Tucker L, Steinberg R, et al.. The Effect of Preservative
331 and Temperature on the Analysis of Circulating Tumor DNA. *Clin Cancer Res Off J Am Assoc Cancer*
332 *Res.* 2017; doi: 10.1158/1078-0432.CCR-16-1691.

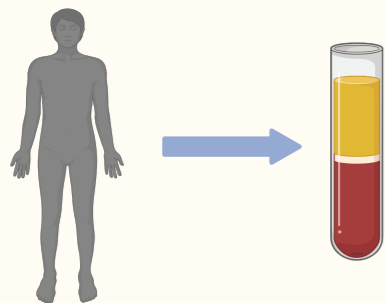
- 333 16. van Dessel LF, Beije N, Helmijr JCA, Vitale SR, Kraan J, Look MP, et al.. Application of circulating
334 tumor DNA in prospective clinical oncology trials – standardization of preanalytical conditions. *Mol*
335 *Oncol.* 2017; doi: 10.1002/1878-0261.12037.
- 336 17. Abbosh C, Birkbak NJ, Wilson GA, Jamal-Hanjani M, Constantin T, Salari R, et al.. Phylogenetic
337 ctDNA analysis depicts early stage lung cancer evolution. *Nature.* 2017; doi: 10.1038/nature22364.
- 338 18. Benjamini Y, Speed TP. Summarizing and correcting the GC content bias in high-throughput
339 sequencing. *Nucleic Acids Res.* 2012; doi: 10.1093/nar/gks001.
- 340 19. Kim CS, Mohan S, Ayub M, Rothwell DG, Dive C, Brady G, et al.. In silico error correction improves
341 cfDNA mutation calling. *Bioinformatics.* 2019; doi: 10.1093/bioinformatics/bty1004.
- 342 20. Esfahani MS, Hamilton EG, Mehrmohamadi M, Nabet BY, Alig SK, King DA, et al.. Inferring gene
343 expression from cell-free DNA fragmentation profiles. *Nat Biotechnol.* Nature Publishing Group;
344 2022; doi: 10.1038/s41587-022-01222-4.
- 345 21. Doebley A-L, Ko M, Liao H, Cruikshank AE, Santos K, Kikawa C, et al.. A framework for clinical
346 cancer subtyping from nucleosome profiling of cell-free DNA. *Nat Commun.* 2022; doi:
347 10.1038/s41467-022-35076-w.
- 348 22. Mathios D, Johansen JS, Cristiano S, Medina JE, Phallen J, Larsen KR, et al.. Detection and
349 characterization of lung cancer using cell-free DNA fragmentomes. *Nat Commun.* 2021; doi:
350 10.1038/s41467-021-24994-w.
- 351 23. Markus H, Contente-Cuomo T, Farooq M, Liang WS, Borad MJ, Sivakumar S, et al.. Evaluation of
352 pre-analytical factors affecting plasma DNA analysis. *Sci Rep.* 2018; doi: 10.1038/s41598-018-25810-
353 0.
- 354 24. Zheng H, Zhu MS, Liu Y. FinaleDB: a browser and database of cell-free DNA fragmentation
355 patterns. *Bioinformatics.* Oxford University Press; 2021; doi: 10.1093/bioinformatics/btaa999.
- 356 25. Zhang W, Wei L, Huang J, Zhong B, Li J, Xu H, et al.. cfDNApipe: a comprehensive quality control
357 and analysis pipeline for cell-free DNA high-throughput sequencing data. *Bioinformatics.* 2021; doi:
358 10.1093/bioinformatics/btab413.
- 359 26. Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, et al.. Sustainable data
360 analysis with Snakemake. F1000Research;
- 361 27. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al.. Twelve years of SAMtools
362 and BCFtools. *GigaScience.* 2021; doi: 10.1093/gigascience/giab008.
- 363 28. Gaspar JM. NGmerge: merging paired-end reads via novel empirically-derived models of
364 sequencing errors. *BMC Bioinformatics.* 2018; doi: 10.1186/s12859-018-2579-2.
- 365 29. Vasimuddin Md, Misra S, Li H, Aluru S. Efficient Architecture-Aware Acceleration of BWA-MEM
366 for Multicore Systems. *2019 IEEE Int Parallel Distrib Process Symp IPDPS.*
- 367 30. Andrews S. FASTQC. A quality control tool for high throughput sequence data.
- 368 31. Pedersen BS, Quinlan AR. Mosdepth: quick coverage calculation for genomes and exomes.
369 *Bioinformatics.* 2018; doi: 10.1093/bioinformatics/btx699.

- 370 32. Ewels P, Magnusson M, Lundin S, Källér M. MultiQC: summarize analysis results for multiple tools
371 and samples in a single report. *Bioinformatics*. 2016; doi: 10.1093/bioinformatics/btw354.
- 372 33. Zohren F, Souroullas GP, Luo M, Gerdemann U, Imperato MR, Wilson NK, et al.. The transcription
373 factor Lyl-1 regulates lymphoid specification and the maintenance of early T lineage progenitors. *Nat*
374 *Immunol*. Nature Publishing Group; 2012; doi: 10.1038/ni.2365.
- 375 34. Jacobs J, Atkins M, Davie K, Imrichova H, Romanelli L, Christiaens V, et al.. The transcription factor
376 Grainy head primes epithelial enhancers for spatiotemporal activation by displacing nucleosomes.
377 *Nat Genet*. Nature Publishing Group; 2018; doi: 10.1038/s41588-018-0140-x.
- 378 35. Chen AF, Liu AJ, Krishnakumar R, Freimer JW, DeVeale B, Belloch R. GRHL2-Dependent Enhancer
379 Switching Maintains a Pluripotent Stem Cell Transcriptional Subnetwork after Exit from Naive
380 Pluripotency. *Cell Stem Cell*. 2018; doi: 10.1016/j.stem.2018.06.005.
- 381 36. Cocce KJ, Jasper JS, Desautels TK, Everett L, Wardell S, Westerling T, et al.. The Lineage
382 Determining Factor GRHL2 Collaborates with FOXA1 to Establish a Targetable Pathway in Endocrine
383 Therapy-Resistant Breast Cancer. *Cell Rep*. 2019; doi: 10.1016/j.celrep.2019.09.032.
- 384 37. Paltoglou S, Das R, Townley SL, Hickey TE, Tarulli GA, Coutinho I, et al.. Novel Androgen Receptor
385 Coregulator GRHL2 Exerts Both Oncogenic and Antimetastatic Functions in Prostate Cancer. *Cancer*
386 *Res*. 2017; doi: 10.1158/0008-5472.CAN-16-1616.
- 387 38. Riethdorf S, Frey S, Santjer S, Stoupiec M, Otto B, Riethdorf L, et al.. Diverse expression patterns
388 of the EMT suppressor grainyhead-like 2 (GRHL2) in normal and tumour tissues. *Int J Cancer*. 2016;
389 doi: 10.1002/ijc.29841.
- 390 39. Reese RM, Harrison MM, Alarid ET. Grainyhead-like Protein 2: The Emerging Role in Hormone-
391 Dependent Cancers and Epigenetics. *Endocrinology*. 2019; doi: 10.1210/en.2019-00213.
- 392 40. Kwan EM, Fettke H, Crumbaker M, Docanto MM, To SQ, Bukczynska P, et al.. Whole blood GRHL2
393 expression as a prognostic biomarker in metastatic hormone-sensitive and castration-resistant
394 prostate cancer. *Transl Androl Urol*. AME Publishing Company; 2021; doi: 10.21037/tau-20-1444.
- 395 41. Kumegawa K, Takahashi Y, Saeki S, Yang L, Nakadai T, Osako T, et al.. GRHL2 motif is associated
396 with intratumor heterogeneity of cis-regulatory elements in luminal breast cancer. *Npj Breast Cancer*.
397 Nature Publishing Group; 2022; doi: 10.1038/s41523-022-00438-6.

398

Data generation

Cell-free DNA from liquid biopsy

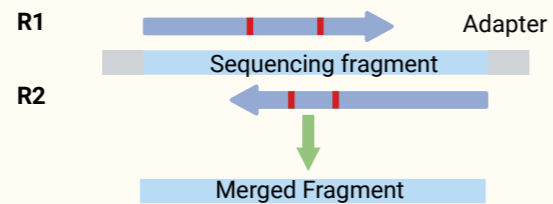


Sequencing



Data pre-processing

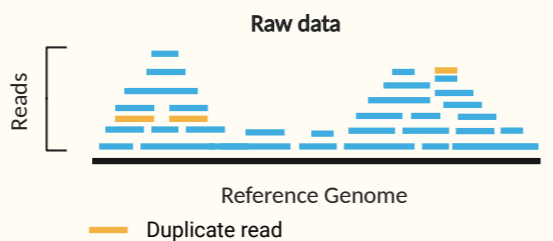
Sequencing error correction and adapter removal



Read length filtering



Mapping and duplicate marking

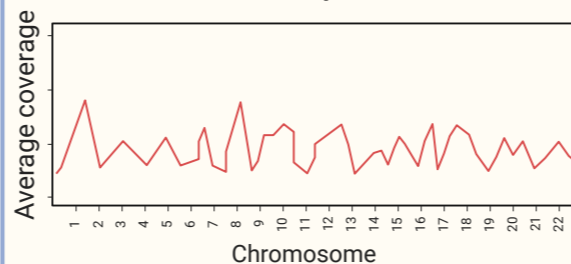


Quality Control

FastQC



Mosdepth

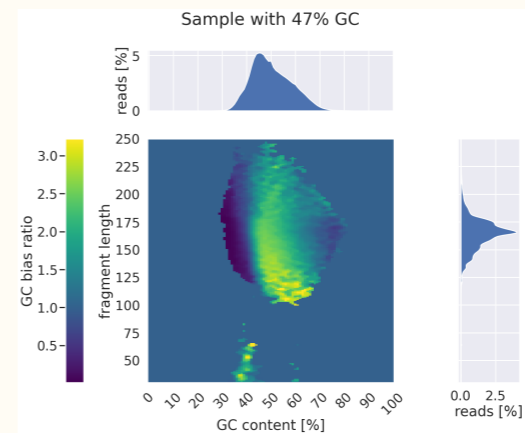


Samtools stats

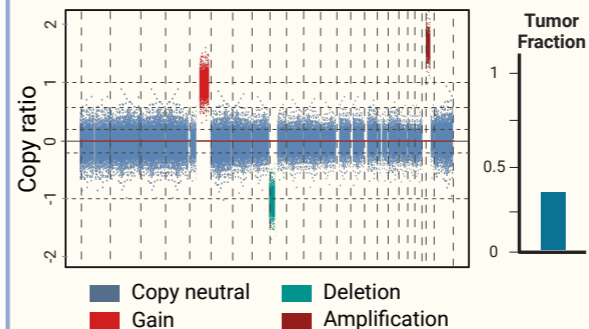
	Total reads	reads MQ 0	Mapped & paired	Duplicated
Sample 1	1040 M	105 M	769 M	130 M
Sample 2	57 M	5 M	42 M	3 M
Sample 3	1275 M	130 M	1052 M	180 M

Sequence analysis

Fragment length based GC bias correction

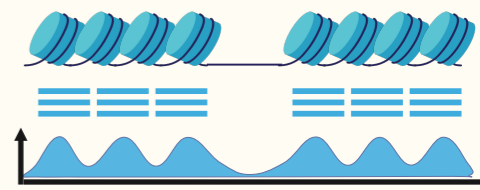


Copy number alteration and tumor fraction estimation (ichorCNA)

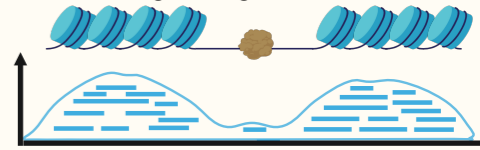


Fragmentomic analysis

Nucleosome footprint

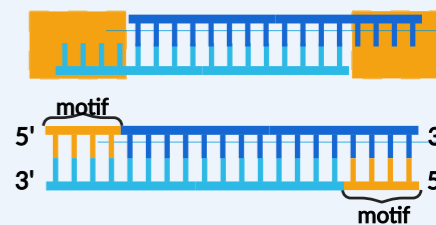


Coverage at regions of interest

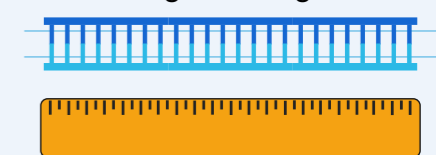


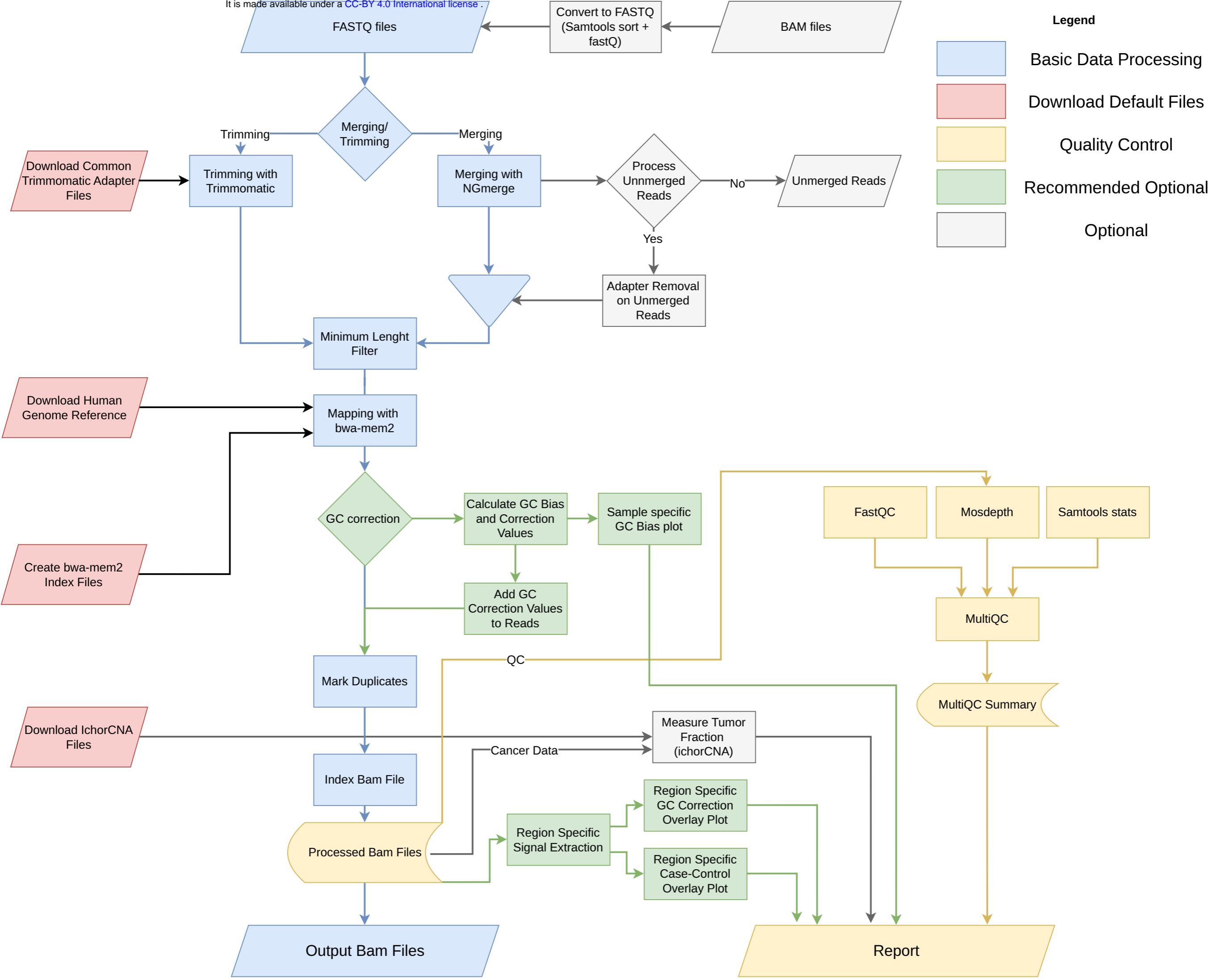
Third-party downstream analyses

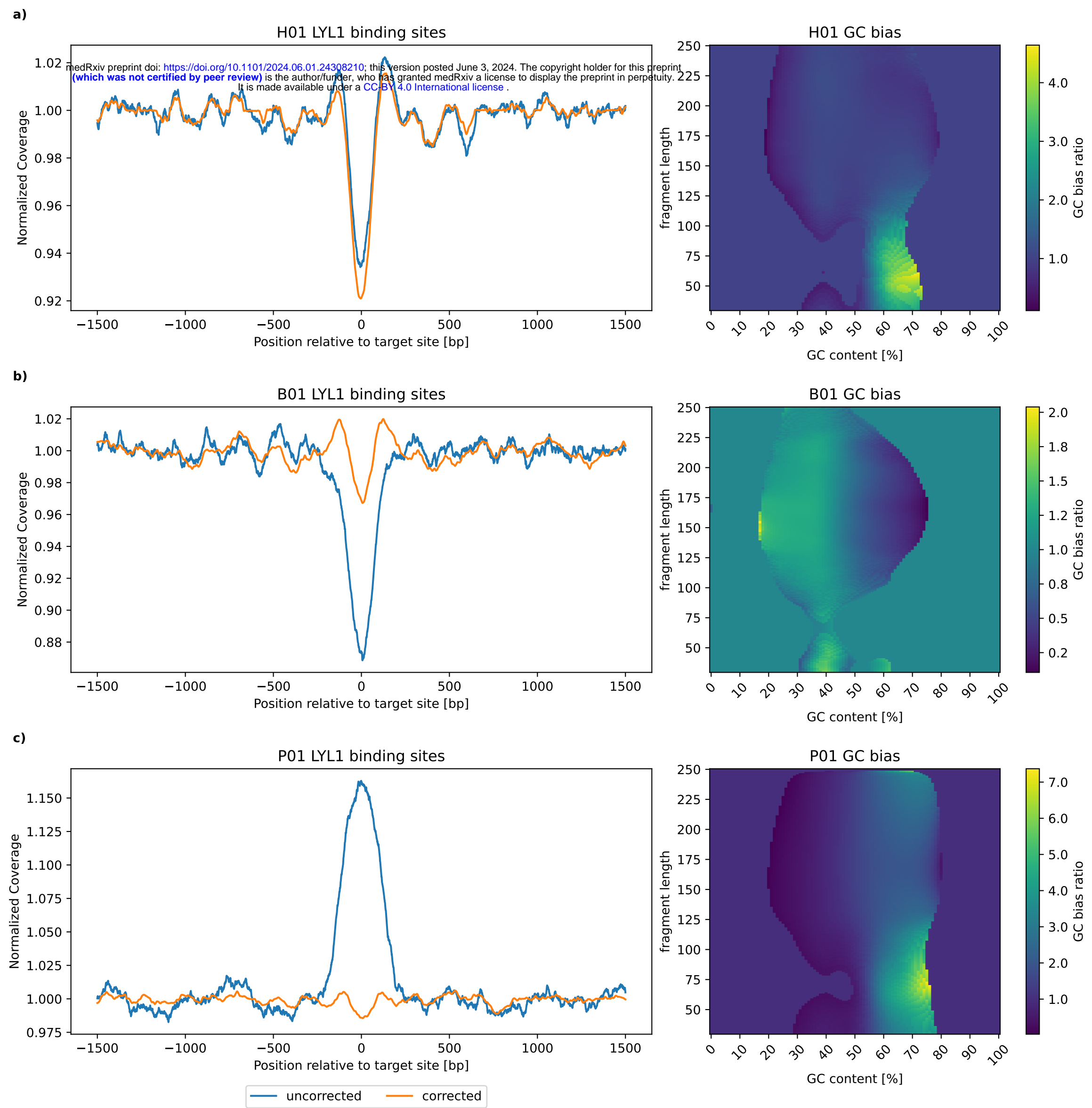
Read-end information

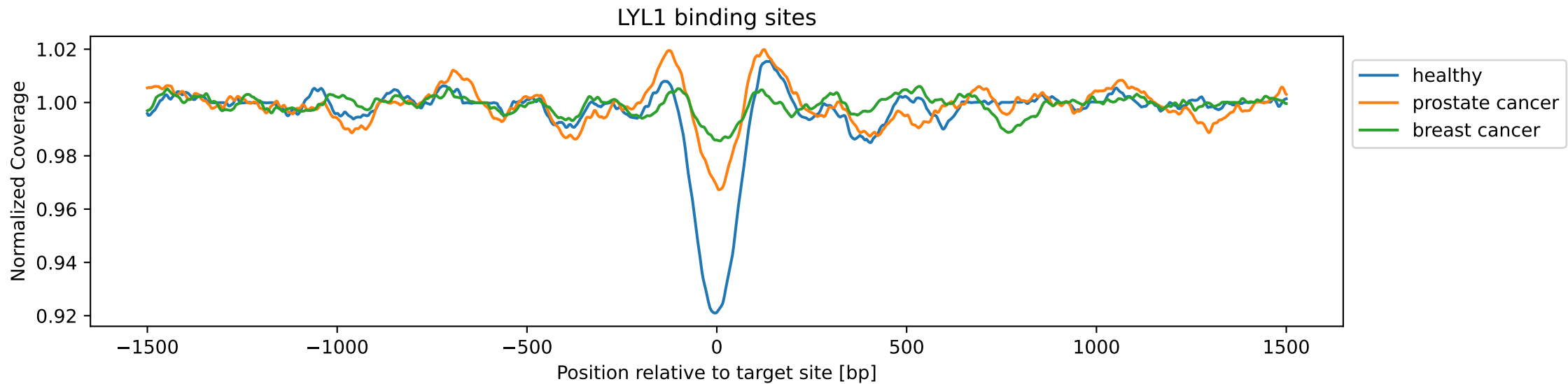


Fragment length







a)**b)**