

1 **Identification and annotation of centromeric hypomethylated regions with**
2 **Centromere Dip Region (CDR)-Finder**

3
4 F. Kumara Mastrorosa¹, Keisuke K. Oshima², Allison N. Rozanski¹, William T. Harvey¹, Evan E. Eichler^{1,3},
5 Glennis A. Logsdon^{1,4,*}

6
7 ¹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA

8 ²Department of Genetics, Epigenetics Institute, Perelman School of Medicine, University of
9 Pennsylvania, Philadelphia, PA, USA

10 ³Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA

11 ⁴Present address: Department of Genetics, Epigenetics Institute, Perelman School of Medicine,
12 University of Pennsylvania, Philadelphia, PA, USA

13
14 ***Correspondence to:**

15 Glennis A. Logsdon, Ph.D.

16 Department of Genetics

17 Epigenetics Institute

18 Perelman School of Medicine at the University of Pennsylvania

19 3400 Civic Center Blvd.

20 Philadelphia, PA 19104

21 Email: glogsdon@penncmedicine.upenn.edu

22 **ABSTRACT**

23 Centromeres are chromosomal regions historically understudied with sequencing technologies due to
24 their repetitive nature and short-read mapping limitations. However, recent improvements in long-read
25 sequencing allowed for the investigation of complex regions of the genome at the sequence and
26 epigenetic levels. Here, we present Centromere Dip Region (CDR)-Finder: a tool to identify regions of
27 hypomethylation within the centromeres of high-quality, contiguous genome assemblies. These regions
28 are typically associated with a unique type of chromatin containing the histone H3 variant CENP-A,
29 which marks the location of the kinetochore. CDR-Finder identifies the CDRs in large and short
30 centromeres and generates a BED file indicating the location of the CDRs within the centromere. It also
31 outputs a plot for visualization, validation, and downstream analysis. CDR-Finder is available at
32 <https://github.com/EichlerLab/CDR-Finder>.

34 **INTRODUCTION**

35 Centromeres are chromosomal regions essential for the segregation of sister chromatids during cell
36 division. Centromeres are typically composed of near-identical tandem repeats known as α -satellite,
37 with other types of satellites (e.g. β -satellite, γ -satellite, HSat1A, HSat1B, HSat2, and HSat3) found
38 within pericentromeric regions^{1,2}. Within the centromere, α -satellite repeats are organized into higher-
39 order repeat (HOR) arrays, which vary in composition and length and constitute the functional unit of
40 the centromeres.

41
42 Recently, long-read sequencing technologies such as Pacific Biosciences (PacBio) high-fidelity (HiFi)
43 and Oxford Nanopore Technologies (ONT) long-read sequencing, as well as newly developed genome
44 assembly algorithms^{3,4}, have led to the reconstruction of the most complex regions of the human
45 genome, including centromeres⁵⁻⁸, telomeres⁹, segmental duplications¹⁰, and other repetitive
46 regions^{9,11}. These⁵⁻⁹ and other¹² studies have also revealed the genetic and epigenetic features of
47 human centromeres. For instance, centromeres are typically hypermethylated throughout the α -satellite
48 HOR array except for a small hypomethylated region called the centromere dip region (CDR)¹². The
49 CDR was shown to coincide with a unique type of chromatin containing the centromeric histone H3
50 variant CENP-A^{6,7}, which marks the site of the kinetochore. This finding was confirmed with several
51 experimental assays⁵⁻⁷.

52
53 Here, we present Centromere Dip Region (CDR)-Finder, a tool to identify and annotate the CDRs
54 based on the CpG methylation data obtained from either PacBio HiFi or ONT data. CDR-Finder detects
55 regions of hypomethylation within sequence-resolved centromeric α -satellite HOR arrays, outputs their
56 coordinates and size, and adds annotations based on their characteristics. It can be used to analyze
57 both large and small centromeres, as long as the centromere is fully assembled and sequence-
58 resolved.

60 **MATERIALS AND METHODS**

61 CDR-Finder requires three input files: a FASTA file of the genome assembly, a BED file containing the
62 coordinates of the region of interest (e.g. the centromere) within the assembly, and a BAM file
63 containing alignments of PacBio HiFi or ONT data with methylation tags to the same genome
64 assembly. CDR-Finder first converts the BAM file to a bedMethyl file using modkit
65 (<https://github.com/nanoporetech/modkit>), which lists the position and frequency of modified bases for
66 each CpG. Then, it divides the region of interest within the methylBED file into sequential 5-kbp bins
67 and calculates the average CpG methylation frequency for each bin. Bins without an assigned value
68 are excluded. Finally, CDR-Finder runs RepeatMasker¹³ on the region of interest to identify the location
69 of the α -satellite sequences (annotated as "ALR/Alpha"), and it calculates the mean CpG methylation
70 frequency across each bin containing α -satellite.

71
72 To identify the CDR(s), our tool first selects bins with a CpG methylation frequency less than the
73 median frequency of all α -satellite sequences in the region of interest. While this frequency can be
74 specified by the user, in our experience, a frequency of 0.34 identifies most CDRs with high precision

75 and recall (described in the example below). Frequencies greater than 0.34 often fail to detect CDRs
76 with shallow hypomethylation, and frequencies less than 0.34 often miscall CDRs due to variation in
77 sequencing coverage. Then, CDR-Finder further refines the bins with a low CpG methylation frequency
78 to those that also have a minimal dip prominence (defined topographically¹⁴) from the median. In our
79 experience, a minimal dip prominence of 0.30 often removes low-confidence calls when the methylation
80 levels are uniformly low across a subregion. Finally, CDR-Finder evaluates the boundaries of each
81 candidate CDR by calculating the mean CpG methylation frequency and then extending each call
82 boundary to the mean CpG methylation frequency +/- a specific value. In our experience, the mean
83 minus one standard deviation is usually sufficient to capture the entire CDR in each call.

84

85 **USAGE AND EXAMPLES**

86 CDR-Finder is organized as a configurable Snakemake pipeline. To run it, the user should first clone
87 the repository from <https://github.com/EichlerLab/CDR-Finder>. Then, the user should modify the
88 configuration file, `config.yaml`, to specify the sample name, the path to sample's genome assembly, the
89 path to the BED file containing a region(s) of interest, and the path to the BAM file containing the
90 alignment of PacBio HiFi or ONT data with methylation tag to the sample's genome assembly.

91

92 In most cases, default parameters can be maintained, but the sequencing coverage and methylation-
93 calling algorithm will affect the ability of CDR-Finder to detect CDRs accurately. As such, the
94 parameters may need to be adjusted accordingly. While CDR-Finder can run on any sequence
95 containing α -satellite DNA, we recommend that the user run it on an α -satellite HOR array with
96 additional flanking sequence on both sides to ensure that the centromere is completely traversed and to
97 observe the transition in methylation patterns between the centromeric and pericentromeric regions.
98 Since CDR detection is based on α -satellite sequences only, the additional flanking sequence will not
99 affect the results of the analysis. We also recommend that the user verify the accuracy of each call
100 based on the coverage data in the plot generated by CDR-Finder and in the original methyl-reads
101 alignment. The user should exclude calls in regions with low coverage and other potential false
102 positives.

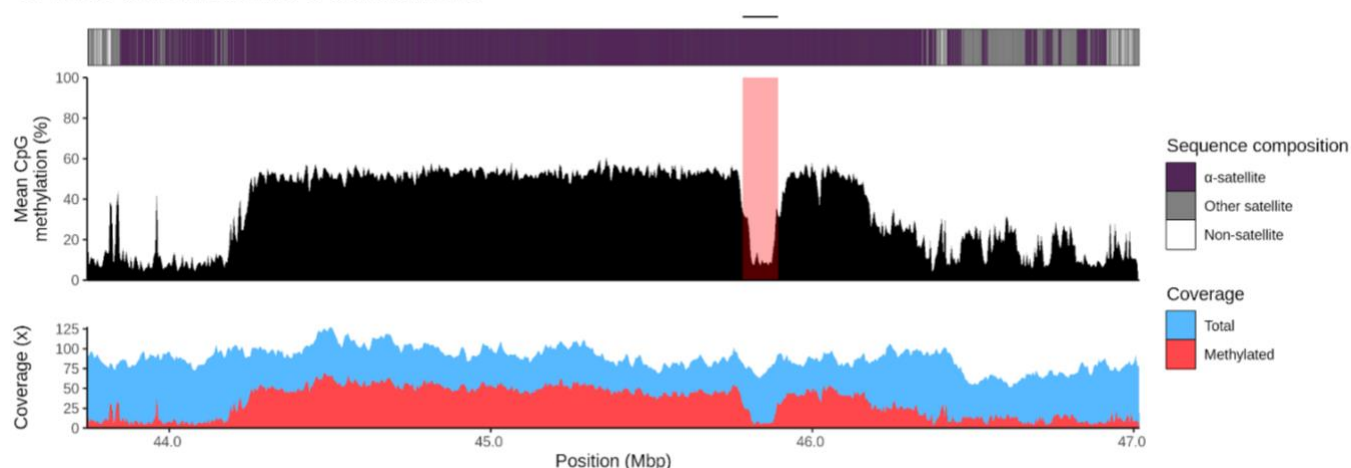
103

104 The pipeline can be run using `conda` or `singularity` with the following commands: `snakemake -np --`
105 `sdm conda -c 4` or `snakemake -np --sdm apptainer conda -c 4`, respectively. It will
106 output a BED file with the coordinates of each CDR call as well as a plot showing the CpG methylation
107 frequency of the region of interest, the overall sequencing coverage and methylated sequencing
108 coverage, and annotation showing the location of the CDR call and the sequence composition of the
109 region (**Figure 1**).

110

111 Below, we show an example of a CDR call for the T2T-CHM13 chromosome 8 centromere⁶
112 characterized by CDR-Finder (**Figure 1**). We used the original chromosome 8 centromere α -satellite
113 HOR array coordinates with ~500 kbp additional sequence on both sides as a target region
114 (`chr8:43,746,447-47,020,471` in the T2T-CHM13 v2.0 genome⁹) as well as the original ONT data
115 generated from the same genome. CDR-Finder called only one CDR, as expected (`chr8:45,789,626-`
116 `45,899,626`, size = 110 kbp). The same T2T-CHM13 chromosome 8 CDR was originally described with
117 a similar size (73 kbp)⁶. However, on that occasion, CpG methylation was detected with a different
118 method, and the CDR was considered as the region with the lowest CpG methylation levels.

CHM13 chromosome 8 centromere



119

120

Figure 1. Detection of a CDR in the T2T-CHM13 chromosome 8 centromere with CDR-Finder. CDR-Finder generates a plot with the annotation of the region (top), mean CpG methylation frequency (middle), and the corresponding read coverage (both total and methylated reads) across the region (bottom). The CDR is highlighted in red and indicated with a black bar on top.

124

To evaluate CDR-Finder's ability to detect CDRs in diverse human centromeres, we ran the tool on 200 completely and accurately assembled centromeres from 15 randomly selected samples from the Human Genome Structural Variation Consortium¹⁵ with default parameters. We manually inspected all calls, counting the number of CDRs that were correctly called, partially called (where the CDR could be extended on one or both sides), incorrectly called, or not called. Our test showed that 98.7% (443/449) of the CDRs were correctly called, indicating a precision of 0.99, and only 43 CDRs were not called (8.7%), indicating a recall of 0.91. In 28 cases (6.3%), the CDR(s) were only partially called, which could be corrected by tweaking the parameters for the specific case. The six erroneous calls could be easily excluded with a visual inspection of the CDR-Finder plots (**Supplementary Table 1**).

134

CONCLUSION

We developed CDR-Finder, a user-friendly method to identify CDRs in complete centromeres. The tool analyzes CpG methylation data from both ONT and PacBio HiFi data and outputs the coordinates of the CDR calls and a plot with related read coverage and highlighted CDR windows. In our experience, both PacBio HiFi and ONT data perform similarly with this tool¹⁶. However, the user can modify the parameters based on specific cases and quality of the data available. We provide an extensive explanation of the parameters with test cases and common issues on the CDR-Finder GitHub page (<https://github.com/EichlerLab/CDR-Finder>).

143

ACKNOWLEDGMENTS

This research was supported, in part, by funding from the National Institutes of Health (NIH) R01HG010169 (E.E.E.) and R00GM147352 (G.A.L.). E.E.E. is an investigator of the Howard Hughes Medical Institute. This article is subject to HHMI's Open Access to Publications policy. HHMI lab heads have previously granted a nonexclusive CC BY 4.0 license to the public and a sublicensable license to HHMI in their research articles. Pursuant to those licenses, the author-accepted manuscript of this article can be made freely available under a CC BY 4.0 license immediately upon publication.

151

CONFLICTS OF INTERESTS STATEMENT

E.E.E. is a scientific advisory board (SAB) member of Variant Bio, Inc. The other authors declare no competing interests.

155

156 **DATA AVAILABILITY**

157 The T2T-CHM13 v2.0 genome assembly, PacBio HiFi data, and ONT data are available at:
158 <https://github.com/marbl/CHM13>. The HGSC genome assemblies, PacBio HiFi data, and ONT data
159 are available at: https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSC3/release. The
160 ONT data used in our tests were basecalled with Guppy v6.3.7, which detects methylated cytosines
161 during basecalling, and were aligned to the T2T-CHM13 v2.0 reference genome using winnowmap2¹⁷.
162 The following command was used for ONT read alignment and processing: `winnomap -W
163 CHM13_repetitive_k15.txt -y --eqx -ax map-ont -s 4000 -t {threads} -I 10g
164 {ref.fasta} {reads.fastq} | samtools view -u -F 2308 - | samtools sort -o
165 {output.bam} -.`
166

167 **REFERENCES**

- 168
- 169 1. Miga, K. H. & Alexandrov, I. A. Variation and Evolution of Human Centromeres: A Field Guide and
170 Perspective. *Annual Review of Genetics* **55**, 583–602 (2021).
- 171 2. Logsdon, G. A. & Eichler, E. E. The Dynamic Structure and Rapid Evolution of Human Centromeric
172 Satellite DNA. *Genes* **14**, 92 (2023).
- 173 3. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly
174 using phased assembly graphs with hifiasm. *Nat Methods* **18**, 170–175 (2021).
- 175 4. Rautiainen, M. *et al.* Verkko: telomere-to-telomere assembly of diploid chromosomes.
176 2022.06.24.497523 Preprint at <https://doi.org/10.1101/2022.06.24.497523> (2022).
- 177 5. Miga, K. H. *et al.* Telomere-to-telomere assembly of a complete human X chromosome. *Nature*
178 **585**, 79–84 (2020).
- 179 6. Logsdon, G. A. *et al.* The structure, function and evolution of a complete human chromosome 8.
180 *Nature* **593**, 101–107 (2021).
- 181 7. Altemose, N. *et al.* Complete genomic and epigenetic maps of human centromeres. *Science* **376**,
182 eabl4178 (2022).
- 183 8. Logsdon, G. A. *et al.* The variation and evolution of complete human centromeres. *Nature* **629**,
184 136–145 (2024).
- 185 9. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
- 186 10. Vollger, M. R. *et al.* Segmental duplications and their variation in a complete human genome.
187 *Science* **376**, eabj6965 (2022).
- 188 11. Rhie, A. *et al.* The complete sequence of a human Y chromosome. *Nature* **621**, 344–354 (2023).
- 189 12. Gershman, A. *et al.* Epigenetic patterns in a complete human genome. *Science* **376**, eabj5089
190 (2022).
- 191 13. Smit, AFA, Hubley, R & Green, P. RepeatMasker Open-4.0. 2013-2015.
- 192 14. Kirmse, A. & de Ferranti, J. Calculating the prominence and isolation of every mountain in the
193 world. *Progress in Physical Geography: Earth and Environment* **41**, 788–802 (2017).

- 194 15. Logsdon, G. A. *et al.* Complex genetic variation in nearly complete human genomes.
195 2024.09.24.614721 Preprint at <https://doi.org/10.1101/2024.09.24.614721> (2024).
- 196 16. Mastroiosa, F. K. *et al.* Complete chromosome 21 centromere sequences from a Down syndrome
197 family reveal size asymmetry and differences in kinetochore attachment. 2024.02.25.581464
198 Preprint at <https://doi.org/10.1101/2024.02.25.581464> (2024).
- 199 17. Jain, C., Rhie, A., Hansen, N. F., Koren, S. & Phillippy, A. M. Long-read mapping to repetitive
200 reference sequences using Winnowmap2. *Nat Methods* **19**, 705–710 (2022).
201