

METHOD

Open Access



HiCAT: a tool for automatic annotation of centromere structure

Shenghan Gao^{1,2†}, Xiaofei Yang^{2,3,4*†}, Hongtao Guo³, Xixi Zhao⁴, Bo Wang^{1,2} and Kai Ye^{1,2,4,5,6*}

[†]Shenghan Gao and Xiaofei Yang contributed equally to this work.

*Correspondence: xfyang@xjtu.edu.cn; kaiye@xjtu.edu.cn

¹School of Automation Science and Engineering, Faculty of Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

²MOE Key Lab for Intelligent Networks & Networks Security, Faculty of Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

³School of Computer Science and Technology, Faculty of Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

⁴Genome Institute, the First Affiliated Hospital of Xi'an Jiaotong University, Xi'an, Shaanxi, China

⁵School of Life Science and Technology, Xi'an Jiaotong University, Xi'an, Shaanxi, China

⁶Faculty of Science, Leiden University, Leiden, The Netherlands

Abstract

Significant improvements in long-read sequencing technologies have unlocked complex genomic areas, such as centromeres, in the genome and introduced the centromere annotation problem. Currently, centromeres are annotated in a semi-manual way. Here, we propose HiCAT, a generalizable automatic centromere annotation tool, based on hierarchical tandem repeat mining to facilitate decoding of centromere architecture. We apply HiCAT to simulated datasets, human CHM13-T2T and gapless *Arabidopsis thaliana* genomes. Our results are generally consistent with previous inferences but also greatly improve annotation continuity and reveal additional fine structures, demonstrating HiCAT's performance and general applicability.

Keywords: HiCAT, Centromere annotation, Long-read sequencing technologies, Gapless genomes

Background

Centromeres play an essential role in the transmission of genetic information between generations. Deep analysis of centromere architecture is critical to understanding genome stability, cell division, and disease development [1]. In most eukaryotes, centromeres exhibit extra-long tandem repeat (TR) sequences, but the sequence and length of repeat units, which are referred to as monomers, vary significantly among species [2]. The canonical order of monomers yields higher-order repeats (HORs) [3]. For example, in the active centromere of the human X chromosome (CENX), 12 monomers (the length of one monomer is approximately 171 bp) are consecutively ordered as HOR units (the length of one HOR unit is approximately 12×171 bp) (Fig. 1a) [4]. The sequence identity between monomers within an HOR unit is only 50–90%, but the pairwise sequence identity between HOR units in a given centromere is as high as 95–100% [5]. The extra-long TRs and high homogeneity make it difficult to achieve accurate assembly of centromeres, hindering thorough investigations of their sequence architecture [5]. The rapid development of long-read sequencing technologies, especially PacBio high-fidelity (HiFi) reads, has greatly improved genome assembly quality [6]. Based on this



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

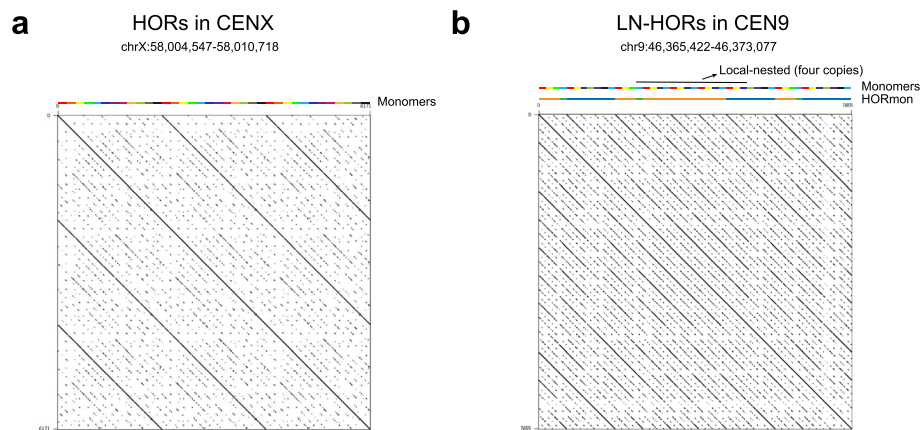


Fig. 1 Examples of higher-order repeats (HORs). **a** HORs in CHM13 CENX. **b** Local-nested HORs (LN-HORs) in CHM13 CEN9. In the monomer tracks, rectangles in various colors represent different monomers. In the HORmon tracks, differently colored rectangles represent different annotations in HORmon. Blue, orange, and green rectangles represent the annotated canonical HORs, partial HORs, and monomers not belonging to any HORs, respectively

progress, the telomere-to-telomere (T2T) consortium presented the complete sequence of the human complete hydatidiform mole (CHM) cell line CHM13 in 2022 [7]. In addition, gap-free genome assembly has been achieved in a few plant genomes, such as those of *Arabidopsis thaliana* and *Oryza sativa* [8, 9]. Significant improvements in genome quality have also contributed to the development of bioinformatic methods for the study of centromere architecture.

Centromere annotation, including monomer inference and HOR detection, is a prerequisite for studying the structure and evolution of centromeres within and between species [10]. Previous studies annotated a substantial number of monomers and HORs in the human genome in a semi-manual manner, facilitating the understanding of centromere architecture [11–13]. However, this semi-manual method lacks a rigorous algorithm definition and is time-consuming and laborious, prohibiting its ready application to new assemblies. To address this question, Dvorkina et al. proposed the first automatic centromere annotation tool, CentromereArchitect [10], which was based on StringDecomposer (SD) [4], an algorithm for detecting sequence blocks by taking monomer templates to decompose centromere DNA sequences. In CentromereArchitect, monomer inference and HOR detection were considered two separate problems without interconnections, which often led to biologically inadequate annotation [14]. The authors next proposed HORmon [14] based on the centromere evolution postulate (CE postulate, where each monomer appears only once in the HOR unit) to address the lack of interconnection issue in CentromereArchitect. HORmon first constructs a *de* Bruijn graph based on monomers inferred from CentromereArchitect and then refines the monomers by considering positional similarity to amend the graph as a single cycle (referred to as the detected HOR) to comply with the CE postulate. Finally, HORmon classifies the detected HORs into canonical and partial HORs. However, the CE postulate has never been strictly proven and heavily depends on parameters [14], while a single occurrence of each monomer in a HOR does not always hold. For example, tandem

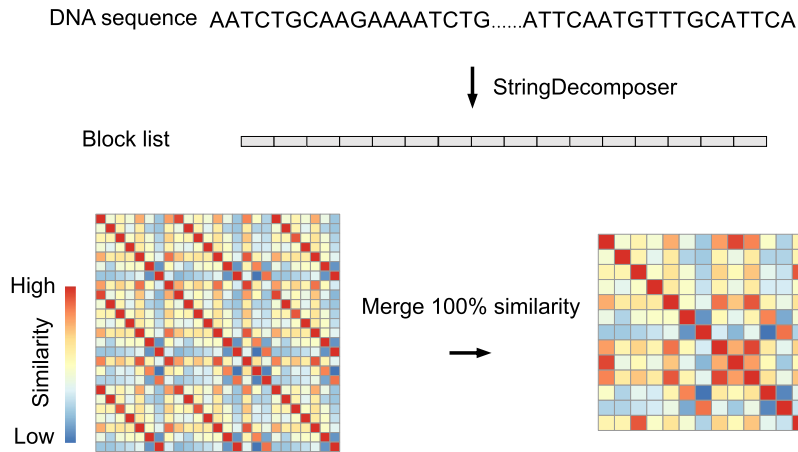
amplification of HOR subunit does occur within HORs and forms so-called local-nested HORs (LN-HORs) (Fig. 1b). Specifically, we observed LN-HORs in CHM13 CEN9, CEN13, and CEN18 (Additional file 1: Fig. S1), violating the CE postulate [14]. Thus, a substantial number of partial HORs were introduced based on the CE postulate, breaking annotation continuity and hindering the characterization of fine internal architectures in these centromeres (Fig. 1b). To overcome these problems, we propose a generalizable automatic centromere annotation tool named HiCAT based on hierarchical tandem repeat mining (HTRM) using a bottom-up iterative TR compression strategy to detect and represent LN-HORs, achieving Hierarchical Centromere structure AnnoTation. We compared HiCAT with HORmon on simulated datasets to validate the performance. And then, we applied HiCAT to newly assembled telomere to telomere (T2T) genomes of human [11] and *Arabidopsis thaliana* [8]. We compared the results from HiCAT and those from semi-manual and HORmon approaches. We found that our automated results are generally consistent with those of previous studies. In addition, HiCAT greatly improved annotation continuity and was able to detect fine structures and organization of HOR units with length variation (LN-HORs) that were missed by other methods. All the comparison results demonstrate the superior performance and generalization of HiCAT.

Results

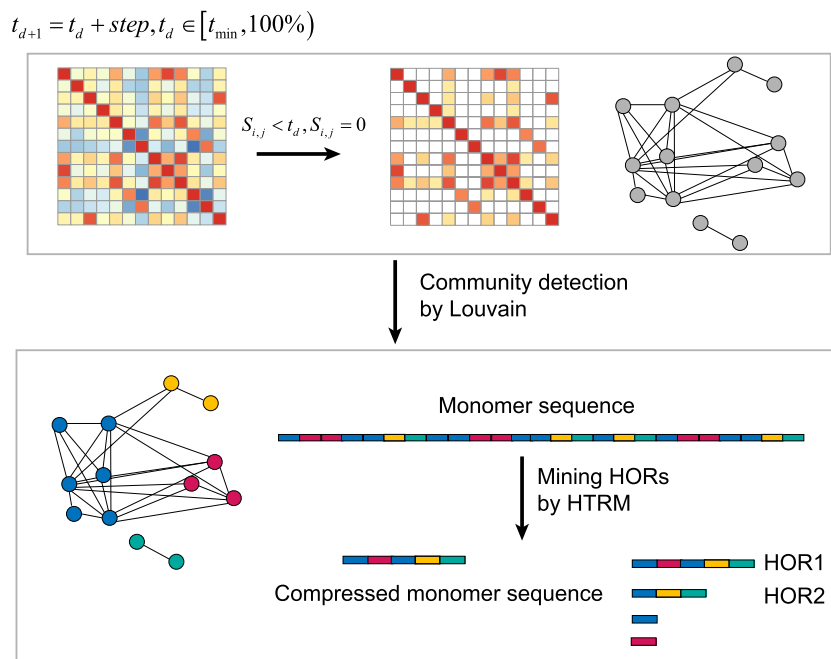
Overview of HiCAT

HiCAT takes a monomer template and a centromere DNA sequence as inputs. There are two steps in HiCAT: generation of a block list and similarity matrix (Fig. 2a) and mining of HORs (Fig. 2b). In the first step, HiCAT uses StringDecomposer [4] to transform a centromere DNA sequence into a block list based on an input monomer template. Each block is a subsequence of the centromere DNA sequence and exhibits high similarity to the monomer template. Then, we defined a similarity score based on the block edit distance to obtain a block similarity matrix (“Methods”). To improve calculation efficiency, we pre-processed the block similarity matrix by merging identical blocks. In the second step, we defined a block graph whose nodes are blocks and edges are links between any two blocks if their similarity value is greater than a given similarity threshold. A series of graphs are created when the similarity threshold iteratively increases from the minimum value (by default 94%) to nearly 100% with a specific step (by default 0.5%). For each constructed block graph, we used the Louvain algorithm [15, 16] to detect block communities, i.e. so-called monomers. We assigned a unique number to each detected monomer as its ID and transformed the block list into a monomer sequence. To detect LN-HORs, we proposed the hierarchical tandem repeat mining (HTRM) method (“Methods”, Additional file 1: Fig. S2 and Additional file 1: Supplementary method). HTRM recursively detected and compressed local TRs in the monomer sequence until no TRs were identified. After HTRM, we merged all TRs with shifted monomer pattern units, such as 1–2–3–4, 4–1–2–3, 3–4–1–2, and 2–3–4–1, to obtain HORs. To build interconnections between monomer inference and HOR detection, we calculated the associated HOR coverage of each similarity threshold and chose the threshold with the largest coverage to obtain HiCAT HORs. Finally, we scored HORs based on coverage

a Generation of a block list and similarity matrix



b Mining higher order repeats



Select the similarity threshold leading to the largest tandem repeat coverage

Fig. 2 Overview of HiCAT. **a** Generation of the block list and similarity matrix. **b** Mining of higher-order repeats (HORs). t_d represents the similarity threshold in the current iteration. t_{d+1} represents the similarity threshold in the next iteration. t_{min} is the minimum similarity threshold. $step$ is the threshold increase for each iteration. For example, $t_0 = 94\%$ and $t_1 = t_0 + 0.5\% = 94.5\%$. $S_{i,j}$ is the similarity between block i and block j . HTRM: hierarchical tandem repeat mining. Colored rectangles in the monomer sequence represent monomers

and the degree of local nesting to rank all HORs (“Methods”). Each HOR was named “R + (rank) + L + (length of HOR unit in the monomer pattern)”. For example, the first HOR in human CENX with 12 monomers was named R1L12.

Performance evaluation on simulated datasets

To comprehensively evaluate HiCAT's performance, we simulated two datasets including both canonical HORs and LN-HORs. The LN-HORs were simulated based on canonical HORs ("Methods", Additional file 1: Fig. S3). We simulated raw monomer template DNA sequence with 100, 200, 300, and 400 bp according to the reported satellite monomers often being 100–400 bp [17]. We built five diverged monomer templates by randomly mutating 10, 20 or 30% bases for each raw monomer template. We generated final DNA sequences by randomly mutating 0.5, 1.5 or 2.5% bases for each diverged monomer template (Additional file 1: Fig. S3). We randomly simulated each case 10 times and obtained 720 simulated HORs and LN-HORs in total ("Methods", Additional file 2: Table S1 and Additional file 2: Table S2).

We first evaluated and compared HiCAT with HORmon (using default parameter) based on simulated canonical HORs. We found that the number of HORs annotated by HiCAT is closer to the ground truth of 40 HOR units than that of HORmon regardless of monomer size, monomer divergence and HOR divergence (Fig. 3a–c and Additional file 2: Table S1). Specifically, HiCAT correctly annotated 40 HOR units in 96.11% (346/360) simulated canonical HORs while 61.67% (222/360) for HORmon (Additional file 2: Table S1). We investigated the special cases that HORmon failed and found that HORmon misannotated majority of the events (89 out of 90 as 33 HOR units) while HiCAT correctly annotated most of the events (82 out of 90 as 40 HOR units) when monomer size was 400 bp, regardless of monomer divergence and HOR divergence (Fig. 3a and Additional file 2: Table S1). Next, we evaluated HiCAT performance on simulated LN-HORs. On average, HiCAT annotated 19.7 local-nested units and 19.9 canonical units (Fig. 3d). Specifically, 90.56% (326/360) LN-HORs from HiCAT annotation were identical to the ground truth of 20 canonical HOR units and 20 LN-HOR units (Additional file 3: Table S2). We further investigated the annotation results of HORmon on the simulated LN-HORs and found that HORmon annotated 11.1 single monomers, 38.5 partial HOR units and 39.3 canonical HOR units for each simulated LN-HOR dataset on average (Fig. 3d).

Overall performance for human CHM13 centromeres

We first applied HiCAT in an active alpha satellite array for each centromere (Additional file 4: Table S3) of the human CHM13-T2T genome (v1.0) [11] and compared the results

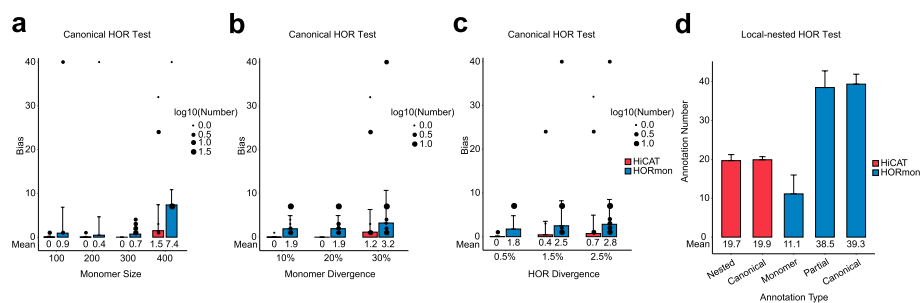


Fig. 3 Simulated validation for HiCAT and HORmon. The bias of canonical HOR annotation of HiCAT and HORmon grouped by **a** monomer size, **b** monomer divergence and **c** HOR divergence. **d** The result of HiCAT and HORmon in local-nested HOR annotation

with published results obtained with semi-manual inference [11, 13]. We found that the HiCAT results were highly consistent with those of previous studies. The reported HORs in 21 out of 23 centromeres (CEN1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 18, 19, 20, 21, 22, and X) were well detected by HiCAT, while substantial differences were observed for the remaining two chromosomes, CEN5 and CEN17 (Additional file 5: Table S4). We first took CEN11 and 15 as examples to explore the HiCAT results. We found that HORs in CEN11 were rather homogeneous, with as few as 12 nested units in R1L5 (Fig. 4a, b). In CEN15, there were approximately four times as many nested units in R1L11 as canonical units (Fig. 4c, d). The monomer pattern of the CEN15 R1L11 unit was 1–2–3–4–5–(6–7–8–9) \times n –10–11. Each number represents a monomer, and “ $\times n$ ” represents the number of times that a defined monomer set was repeated. For example, four consecutive monomers 6–7–8–9 in the R1L11 unit repeated, and most of them have two copies, while other numbers of repeats also existed (Fig. 4e).

In CEN1, 8, 9, 10, 13, and 19, previously reported HORs were not ranked first but among the top five HiCAT results (Additional file 5: Table S4) due to local-nested tandem amplification (Fig. 4 and Additional file 1: Fig. S4). For CEN1, the first HiCAT HOR was R1L2 with two monomers, which was consistent with previously reported dimeric expansions in D1Z7 [13] (the HOR name in previous studies was displayed as “D + chromosome number + Z + sequential number” [3, 11]). In the CHM13 genome, a 1.7-Mb inversion in the CEN1 active alpha satellite array [11] and split the reported D1Z7 into two HORs, R2L6, and R3L6 (Fig. 4f), with reversed monomer patterns of 1–2–5–6–4–3 and 3–4–6–5–2–1, respectively (Fig. 4g). In R2L6 and R3L6, we also detected LN tandem amplification in two monomers (1–2) at a peak of four copies (Fig. 4g).

For CEN8, HiCAT detected four frequent HORs, namely, R1L15, R2L7, R3L8, and R4L11 (Fig. 4h, i), corresponding to the reported HORs with different number of monomers (4, 7, 8, and 11) [11, 18]. We found that R1L15 reported by HiCAT is a combination of two HORs with 11 monomers and 4 monomers, respectively (Additional file 1: Fig. S5a, b). We also detected the reported location bias of these HORs [18], that is R4L11 was mainly distributed in the marginal area, while R2L7 was enriched in the center. R1L15 and R3L8 were distributed between R4L11 and R2L7. Furthermore, we compared the monomer annotation from HiCAT and previous reported D8Z2 (so-called S2C8H1L) [11, 18], and found monomer annotation is largely consistent although S2C8H1L.7 and S2C8H1L.4/7 were annotated as monomer 11 in HiCAT due to their high identity (95.93%) (Additional file 1: Fig. S5a, c). All of these results demonstrate the reliability of HiCAT on centromere HOR annotation. From these HOR patterns detected by HiCAT, we concluded that HOR patterns in human centromeres are diverse and some of them show location specificity. Moreover, the results showed that large structure rearrangements exist in some centromeres.

Substantial differences between HiCAT and semi-manual HOR annotations in CEN5 and CEN17

Previous studies have reported that CEN1, 5, and 19 contain shared HORs with six monomers (D1Z7, D5Z2 and D19Z3, so-called S1C1/5/19H1L) belonging to suprachromosomal family 1 (SF1) and are organized as alternating dimers of J1 and J2 monomers [13]. D1Z7 and D19Z3 were detected in CEN1 (R2L6 and R3L6) and

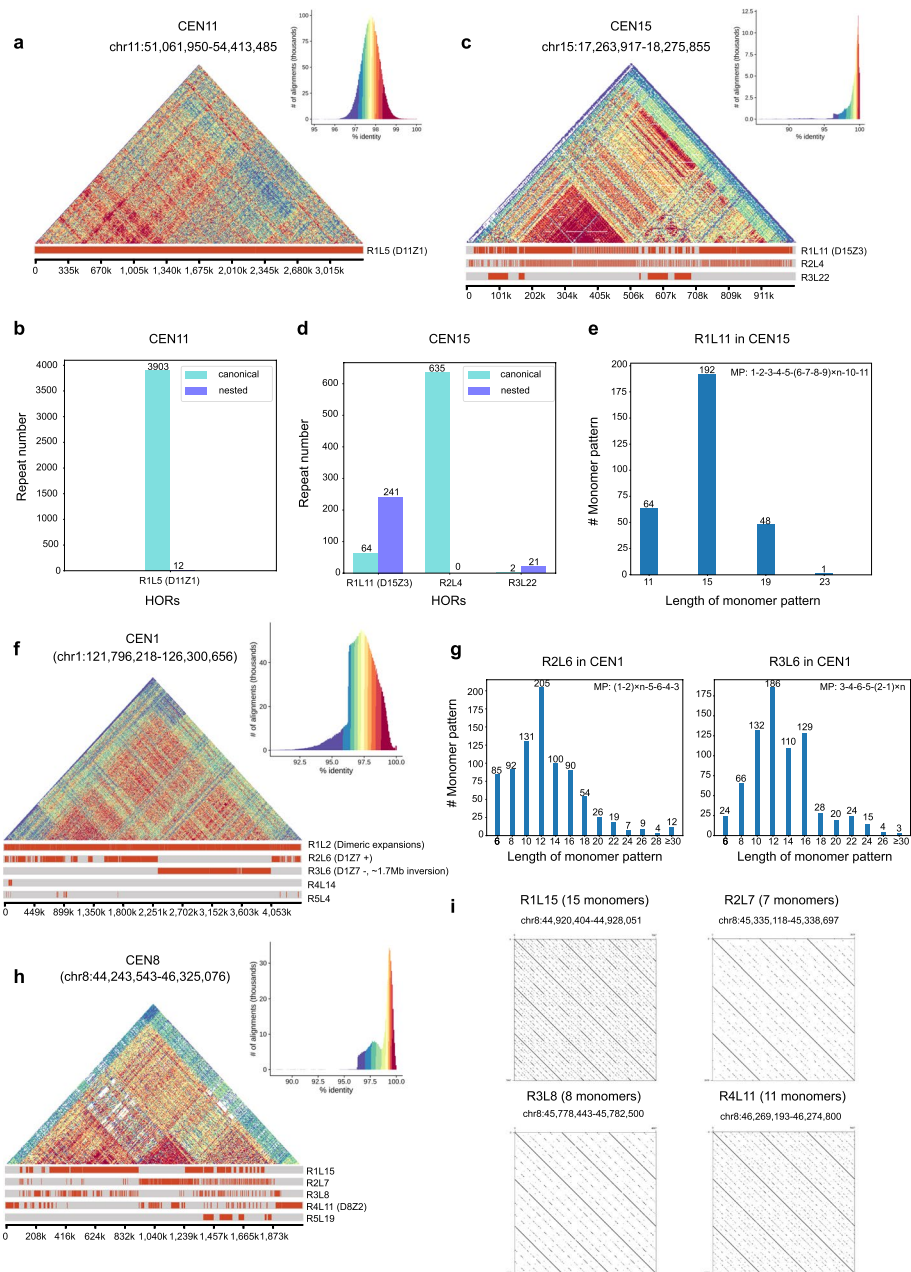


Fig. 4 Fine structures in CHM13 CEN11, CEN15, CEN1, and CEN8. **a** Structure and annotation of CEN11. **b** The numbers of HOR repeats in CEN11. **c** Structure and annotation of CEN15. **d** The numbers of HOR repeats in CEN15. **e** The numbers of monomer patterns in CEN15 R1L11. **f** Structure and annotation of CEN1. **g** The number of monomer patterns in CEN1 R2L6 and R3L6. **h** Structure and annotation of CEN8. **i** Dotplots for different HORs in CEN8. D11Z1, D15Z3, D1Z7, and D8Z2 are previously reported HORs. MP is the monomer pattern. # means the number of

CEN19 (R2L6), respectively (Fig. 4f and Additional file 1: Fig. S4a), while D5Z2 was not detected in the top five HiCAT results in CEN5 (Additional file 1: Fig. S6a). CEN5 contains hybrid monomers which are a concatenate of two or more monomers, and D5Z2 is the result of dehybridization [11, 13, 14]. In HiCAT annotation, the top pattern in CEN5 was R1L12 with 12 monomers, which are consistent with monomers

in D5Z2 annotated from Altemose et al. [11] (Additional file 1: Fig. S6b). HiCAT annotated three major hybrid monomers S1C1/5/19H1L.2/6, S1C1/5/19H1L.6/4 and S1C1/5/19H1L.2/4 as monomer 2, 7, and 8, respectively (Additional file 1: Fig. S6b,c). R1L12 was D5Z2 without dehybridization (Additional file 1: Fig. S6c,d). In R1L12, the number of nested units was approximately three times greater than that of canonical units (Fig. 5a). The HOR patterns with monomer lengths of 12, 16, and 20 were the top-three most frequent types of patterns, and their specific patterns were 1–2–3–4–5–6–1–7–(1–2–3–8) × *n* with *n* = 1, *n* = 2 and *n* = 3, respectively (Fig. 5b, Additional file 1: Fig. S6d). The canonical HORs in R1L12 were significantly enriched (*p*-values < 0.05, *z*-test) in the marginal area and LN-HORs including 1–2–3–4–5–6–1–7–(1–2–3–8) × 2 and 1–2–3–4–5–6–1–7–(1–2–3–8) × 3 were significantly enriched (*p*-values < 0.05, *z*-test) in the center (“Methods”, Fig. 5c and Additional file 6: Table S5).

Two HORs, D17Z1-B and D17Z1, were reported in CEN17 [11, 14]. D17Z1-B with 14 monomers was detected as R3L14 by HiCAT (Fig. 5d), while D17Z1 with 16 monomers was detected as a special case of HiCAT R1L14 with monomer pattern of 1–2–3–4–5–6–(7) × *n*–8–9–10–11–12–13–14 (Fig. 5e). For R1L14, 1272 HOR units contain local-nested TRs, while as few as 10 units were canonical (Fig. 5e, Additional file 1: Fig. S6e). Most of the R1L14 units contain 16 monomers with *n* = 3 (Fig. 5e), corresponding to D17Z1 with monomer W1 (monomer 7 in HiCAT) triplication (Additional file 1: Fig. S6f) [19]. Moreover, we also detected other rare fine structures of R1L14 with different copy numbers of monomer 7 up to five (Additional file 1: Fig. S6g).

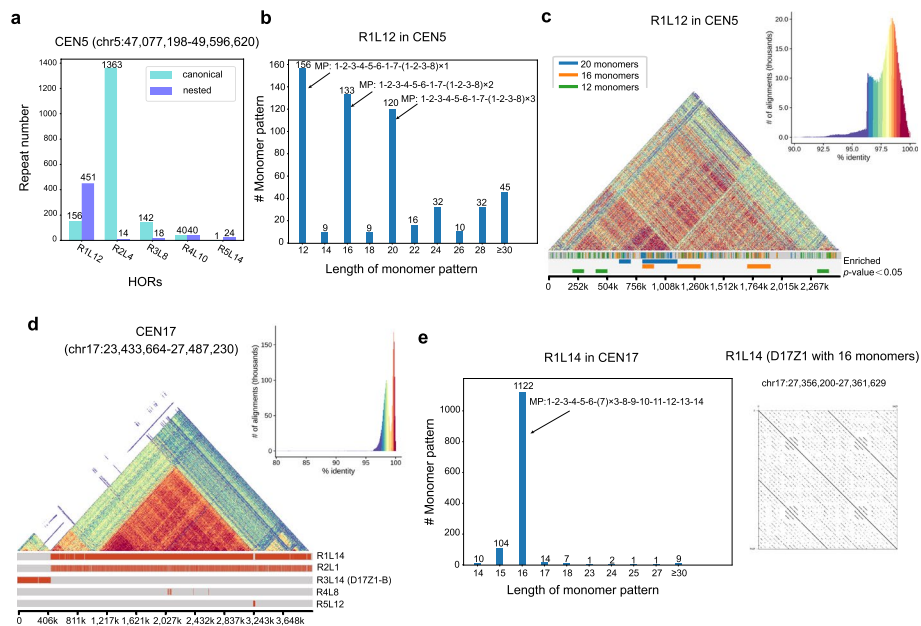


Fig. 5 Resolving centromere structure in CHM13 CEN5 and CEN17. **a** The HOR repeat number in CEN5. **b** The number of monomer patterns in CEN5 R1L12. **c** Structure and annotation of CEN5 for R1L12 with different monomer pattern lengths. We split CEN5 into 25 bins and labelled the significantly enriched (*p*-value < 0.05, *z*-test) R1L12 units for each bin. **d** Structure and annotation of CEN17. **e** The number of monomer patterns in CEN17 R1L14 and dot plot for R1L14 (D17Z1) with 16 monomers. D17Z1 and D17Z1-B are previously reported HORs in CEN17. MP is the monomer pattern. # means the number of

Comparison with HORmon annotation on human CHM13 centromere

We also compared the HORs detected by HiCAT and HORmon [14]. First, we evaluated centromere annotation coverage and continuity in all CHM13 centromeres (Additional file 7: Table S6) and found that the median coverage of both methods was greater than 98% (Fig. 6a). Moreover, we found that HiCAT significantly outperformed HORmon (p -value = $4.6e - 7$, Wilcoxon rank sum test) in terms of continuity, with fewer annotation breakpoints because the LN-HORs were well captured by HTRM (Fig. 6b and Additional file 1: Fig. S7a).

Next, we further compared HiCAT with HORmon in detail by examining CEN9, 13 and 18, which have extensive LN-HORs. Overall, the monomers inferred by the two methods were largely consistent (Fig. 6c and Additional file 1: Fig. S7b, c). For example,

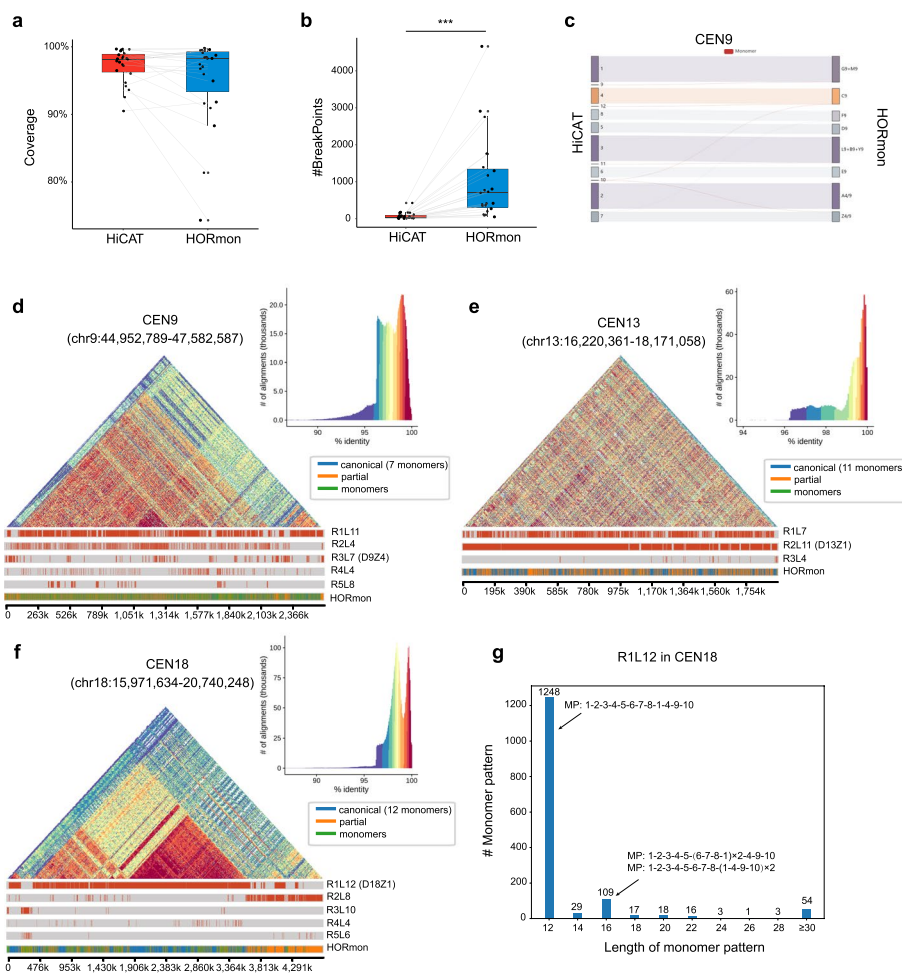


Fig. 6 Comparison of HOR annotations between HiCAT and HORmon. **a** Compared with HORmon in annotation coverage. **b** Compared with HORmon in annotation continuity. p -value = $4.6e - 7$. *** represents p -value < 0.001, Wilcoxon rank sum test. **c** Monomer Sankey plot for CEN9 showing the high consistency between the two methods. To display the frequent monomers, we filtered the links with fewer than 10 matches. The complete Sankey plots are shown in Additional file 1: Fig. S7g-i. **d-f** Structure and annotation of CEN9 (**d**), CEN13 (**e**), and CEN18 (**f**) with two methods. Here, “canonical” represents canonical HORs, “partial” represents partial HORs, and “monomers” represents monomers that do not belong to any HORs. **g** The number of monomer patterns in CEN18 R1L12. D9Z4, D13Z1, and D18Z1 are previously reported HORs. MP is the monomer pattern. # means the number of

the frequent monomers inferred by HiCAT and HORmon were consistent in CEN9 (Fig. 6c) but different in CEN13 monomer 1 and CEN18 monomers 1 and 4 due to a few-nucleotide difference (Additional file 1: Fig. S7b, c) [14].

For HOR detection, HORmon detected canonical HORs with a monomer pattern as A4/9-(L9 + B9 + Y9)-C9-D9-E9-Z4/9-(G9 + M9) in CEN9 [14]. However, monomer F9 with a frequency of 1193 was annotated as a single monomer in HORmon not belonging to any HORs, reducing the coverage of HOR annotation. In HiCAT, due to the HTRM method, monomer 8 (corresponding to monomer F9 in HORmon) was annotated as a subcomponent of R1L11 with a monomer pattern of $(1-2-3-4) \times m-5-6-7-(1-2-3-8) \times n$ (Fig. 6d), resulting in an increase in coverage from 88% (in HORmon) to 94% (in HiCAT) (Additional file 1: Fig. S4e). In CEN13 and CEN18, the monomer patterns of HORs were consistent between HORmon and HiCAT; e.g. D13Z1 (HORmon) equalled R2L11 (HiCAT) in CEN13, and D18Z1 (HORmon) equalled R1L12 (HiCAT) in CEN18 (Fig. 6e, f). However, nearly half of the regions were defined as partial HORs or single monomers by HORmon in CEN13 and CEN18 (Additional file 1: Fig. S7d), generating 726 and 1750 breakpoints, respectively, more than 10 times the number in HiCAT (Additional file 7: Table S6). We reported more fine structures of HORs than HORmon. For example, the canonical monomer pattern R1L12 in CEN18 was 1-2-3-4-5-6-7-8-1-4-9-10, and most of the nested units contained 16 monomers and 6-7-8-1 or 1-4-9-10 in the R1L12 unit repeated (Fig. 6g). Interestingly, we found that the HOR R2L8 in CEN18 with monomer pattern 1-2-3-4-5-6-7-8 was mainly concentrated on the right end of CEN18, reported as partial HORs in the HORmon annotation (Fig. 6e, Additional file 1: Fig. S7f).

Annotation of centromere structures in the plant genome

To demonstrate generalization of HiCAT, we applied it to *Arabidopsis thaliana* Col-CEN centromeres assembled by Naish et al. [8]. We first evaluated the accuracy of HOR annotation by comparing our results with the reported representative HOR region of chr2:4,808,994–4,826,785 [8]. HiCAT detected this HOR as R18L8 (chr2:4,800,609–4,844,007) with a canonical monomer pattern of 2-1-3-4-2-3-4-1 (Fig. 7a, b, Additional file 1: Fig. S8, Additional file 8: Table S7). Next, we applied HiCAT to all centromeres in the Col-CEN assembly (Additional file 4: Table S3, Additional file 8: Table S7). In contrast to human centromeres, in which most HORs evolved from dimers or pentamers [19], we found one monomer tandem amplification (monomic expansion) in all *Arabidopsis thaliana* centromeres (Fig. 7c, Additional file 1: Fig. S9). For example, in CEN1, the top HOR was R1L2 with canonical pattern 1-2 (Fig. 7c, d), and monomers 1 and 2 had a substantial number of copies (Fig. 7e).

Discussion

High-precision and long-read sequencing technologies have revolutionized genome assembly, unlocking complex centromere regions and signalling a new stage in genomics research. The new computing problems introduced by these advances, such as the centromere annotation problem, require novel bioinformatics methods. Here, we propose HiCAT, a generalized computational tool based on the HTRM method to automatically process centromere annotations. The simulated tests demonstrated that HiCAT

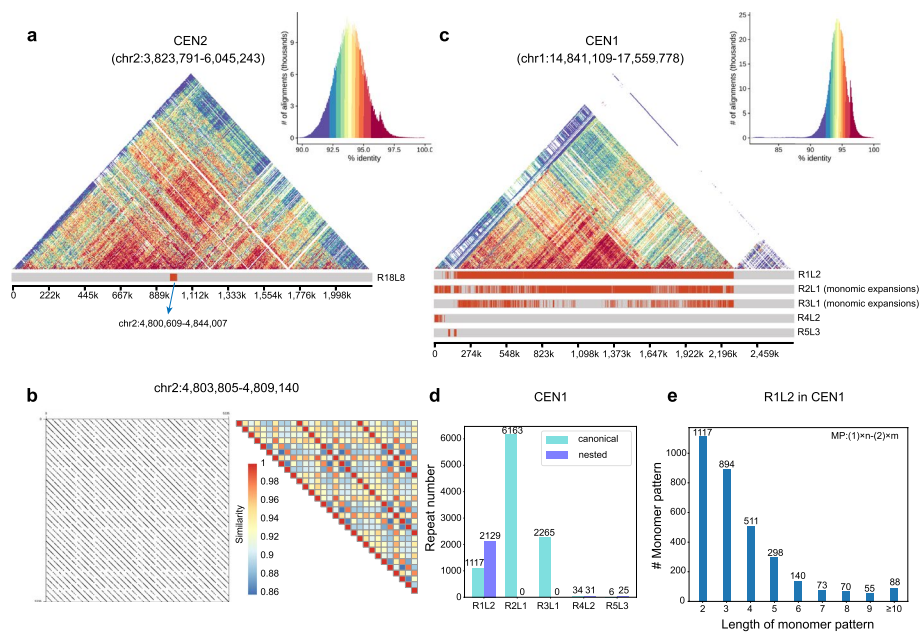


Fig. 7 Annotation of centromere structures in *Arabidopsis thaliana* CEN2 and CEN1. **a** Structure and annotation of CEN2 R18L8. **b** Dot plot and similarity heatmap for a part of R18L8. The complete dot plot and similarity heatmap are shown in Additional file 1: Fig. S8. **c** Structure and annotation of CEN1. **d** The HOR repeat number in CEN1. **e** The number of monomer patterns in CEN1 R1L2. MP is a monomer pattern. # means the number of

outperforms HORmon. Furthermore, HiCAT is able to correctly annotate HORs in both human and plant centromeres, especially those with complex LN-HORs which were annotated as fragmented partial HORs in HORmon. This improvement facilitates the detection of fine structures and organization of HOR units with length variation in centromeres. For example, in human centromere, we found that nested tandem amplification of dimer is common in CEN1 at a peak of four copies (Fig. 4g). In CEN5, we found that LN-HORs were enriched in the center and canonical HORs were significant enriched in the marginal area (Fig. 5c). In CEN17, HiCAT reported monomer 7 with different numbers of copies up to five (Additional file 1: Fig. S6g). HiCAT contributes to the growing but limited toolkit of methods needed to annotate newly assembled satellite DNA regions in both human and non-human genomes. HiCAT results illuminate the organization and evolution of centromere-associated satellite DNA in ways that semi-manual annotation might miss, while enabling the investigation of these regions in an automated fashion in many genomes.

The efficiency of any computational approach is vital for its success. We ran HiCAT on a Linux machine with 28 cores (Intel(R) Xeon(R) Gold 6132 CPU @ 2.60 GHz). In all our tests, the maximum runtime was approximately 2 h for *Arabidopsis thaliana* CEN5, with a length of 2.8 Mb, and the minimum runtime was only 28 s for CHM13 CEN21, with a length of 331 kb (Additional file 9: Table S8).

Currently, HiCAT annotates both canonical and local-nested HORs (LN-HORs) while the latter was defined as tandem amplification of a subunit (Additional file 1: Fig. S10), such as from A-B-C-D-E-F (canonical) to $(A-B) \times n - C-D-E-F$ (local-nested). For the special case of A-B-A-B-C-D-E-F, HiCAT annotates it as

$(A-B) \times 2-C-D-E-F$ while a single deletion of $C-D-E-F$ indeed results in the same outcome (Additional file 1: Fig. S10b, c). We admit that in this case, both interpretations are equivalent. However, we observed abundant cases with $n > 2$, like $A-B-A-B-A-B-C-D-E-F$. For example, in CEN1, two monomers 1–2 in HOR pattern 1–2–5–6–4–3 are tandem amplification and most of them have four copies (Fig. 4g). We also detected other patterns of 1–2–5–6–4–3 with different copy numbers of 1–2 even more than ten (Fig. 4g). These cases would seem more likely to be due to nested tandem repeat amplification rather than by a multiple deletion process (Additional file 1: Fig. S10d, e). Having observed above instances, we believe that a unified framework of tandem repeat amplification to annotate both $n = 2$ and $n > 2$ would be preferred, but we do not rule out alternative explanations, like deletions.

As promising as HiCAT is, there are still some technical limitations and future work that we plan to address. The first is the parameter of minimum similarity. Although we selected the max TR coverage result to guide selection of the similarity threshold, some concerns still should be discussed. If the minimum similarity threshold is set too low, some monomers may be merged, and we may obtain the ancestral state. If the parameter is set too high, the similarity judgment between blocks will be too strict, resulting in too many monomers and leading to failure of HOR detection. In our research, based on a previous study in human centromeres, we set the minimum similarity threshold as 94% since the similarity between HOR units was reported to be in the range of 95–100% in humans [5]. For future newly assembled genomes, this parameter may need to be adjusted to adequately reflect centromere evolution. Although HOR detection from a fully assembled genome gives us comprehensive centromere structures, generating a full genome assembly is still a challenging problem. Annotation of HORs from raw reads is one possible way to obtain and validate centromere structures, and the method named Alpha-CENTAURI has been proposed and applied [20, 21]. We will update HiCAT to accept raw reads as input to extend its application scenarios. Finally, hybrid monomers are also important for comprehensively studying centromere architecture and evolution and were hypothesized as the “birth” of new frequent monomers and were reported in human CEN5 and CEN8 [14]. Currently, HiCAT defines monomers based on only the community detection algorithm, and we will update the monomer inference step to detect hybrid monomers in the future.

Conclusions

We have presented a generalized computational tool, HiCAT based on the HTRM method to automatically process centromere annotations. In human and *Arabidopsis thaliana* centromeres, we showed that HiCAT annotation not only were generally consistent with previous inferences but also greatly improved annotation continuity and revealed additional fine structures, demonstrating HiCAT’s performance and general applicability. We believe that with the emergence of a substantial number of high-quality genomes, HiCAT will promote the study of pan-species centromere diversity and genetic diseases due to defects in centromeres.

Methods

Datasets in humans and *Arabidopsis thaliana*

We obtained active alpha satellite arrays from the complete sequence of the human CHM13 cell line assembled by the T2T Consortium (version 1.0) [11, 14]. CHM13 centromere annotation by the T2T Consortium was obtained from UCSC Table browser. HORmon annotation of CHM13 centromeres was downloaded from figshare [14]. We used the Col-CEN assembly of the *Arabidopsis thaliana* genome and obtained the corresponding centromere coordinates from Naish et al. [8]. The centromere regions in both CHM13 and Col-CEN are summarized in Additional file 4: Table S3.

Generation of the block list and similarity matrix

The first step of HiCAT was to decompose the centromere DNA sequence into the block list based on the input monomer template by StringDecomposer [4] (Fig. 2a). We defined the similarity between blocks b_1 and b_2 as:

$$1 - ed(b_1, b_2) / \max(b_1.len, b_2.len) \quad (1)$$

where ed is edit distance between b_1 and b_2 . $b.len$ is the block length. We calculated the similarity of each block pair to obtain the similarity matrix. Then, we merged the identical blocks (similarity = 100%) to obtain the merged similarity matrix for improving computing efficiency in the HOR mining step.

Mining HORs

Based on the merged similarity matrix, we first defined the block graph, whose nodes are blocks and edges are links between any block pairs if their similarity is greater than a given similarity threshold. A series of block graphs were constructed based on the similarity threshold iteratively increasing from the minimum value (by default 94%) to nearly 100% with a specific step (by default 0.5%). Then, we applied the Louvain algorithm [15, 16] to detect communities in each graph and considered each detected community as a monomer. We assigned a unique number to each monomer as its ID. Next, we transformed the block list into monomer sequences based on block communities (Fig. 2b). Since local-nested TRs hinder the detection of HORs, we proposed the HTRM method to iteratively detect TRs in monomer sequences. HTRM includes monomer tandem detection, region checking, and sequence updating modules. The input of HTRM is a monomer sequence with an upper bound for the length of the TR unit (by default 40 for improving efficiency). We defined a top layer data structure to record non-overlapping TRs with maximum coverage. First, HTRM applied a monomer TR detection module (Additional file 1: Fig. S2a and Additional file 1: Supplementary method) to detect new TRs with a given TR unit length. The initial TR unit length is one. In the second step, we performed region checking (Additional file 1: Fig. S2b) to check for overlap between newly detected TRs (new TRs) and TRs already stored in the top layer (old TRs). The new TRs and old TRs were modified based on four situations. If there was no overlap between them, the new TRs could be saved in the top layer directly. If partial overlap was detected between old and new

TRs, the overlapping new TRs were removed, and the remaining ones were saved in the top layer. If new TRs covered old TRs, the new TRs replaced old TRs in the top layer. Finally, if new TRs were covered by old TRs, the new TRs were discarded. In the sequence updating module, if the top layer was not updated in the region checking step, the TR unit length for detection was increased by one to redetect TRs. Otherwise, the monomer sequences of the newly saved TR region were compressed. After compression, we redetected the TRs by resetting the TR unit length to one. The details and pseudocode of HTRM are shown in the Additional file 1: Supplementary method. After HTRM, all detected TRs are reported, and their units are normalized; e.g. units of 4-1-2-3, 3-4-1-2, and 2-3-4-1 will be normalized as 1-2-3-4. Then, we merged TRs with the same ordered set of normalized units as a HOR. We calculated the associated HOR coverage of each similarity threshold and chose the threshold with the largest coverage for defining HiCAT HORs. Finally, we ranked HiCAT HORs by HOR score combining the coverage and the degree of local nesting. The HOR score is defined as:

$$HORscore = cr * pr \quad (2)$$

$$cr = HOR.len / m.len \quad (3)$$

$$pr = HOR.rn / (HOR.len / HORunit.len) \quad (4)$$

where cr is the coverage for the HOR in the input monomer sequence. pr represents the degree of local nesting. $HOR.len$ is the length of the HOR region in the monomer pattern, and $m.len$ is the length of the monomer sequence. $HOR.rn$ is the repeat number for the HOR, and $HORunit.len$ is the length of the HOR unit in the monomer pattern. If the HOR is over-compressed, which means that it contains only a small number of repeats but with high coverage, $HOR.rn$ will be significantly smaller than $HOR.len / HORunit.len$, and pr will balance the coverage and nested degree of the HOR. We named each HOR in each chromosome as “R + (ranking) + L + ($HORunit.len$)”. For example, in human CEN11, the first HOR is R1L5.

Simulated tests

We simulated two datasets including canonical HORs and LN-HORs. The simulation process contains two steps: simulating monomer sequences and simulating DNA sequences (Additional file 1: Fig. S3a). In the first step, we set HOR unit with five monomers and simulated 40 HOR units as canonical HOR monomer sequences. Next, we randomly selected 20 out of 40 canonical HOR units to generate LN-HORs. For each chosen unit, we randomly selected one to four consecutive monomers of the unit and tandemly amplified the chosen monomer(s) for a random number of times within the range from two to five. In the second step, we simulated raw monomer template DNA sequence with 100, 200, 300, and 400 bp since satellite monomer size are often observed in 100–400 bp [17], respectively. Then, we generated five diverged monomer templates by randomly mutating 10, 20, or 30% bases for each raw monomer template. Finally, based on the monomer sequences acquired in first step, we generated final DNA

sequences by randomly mutating 0.5, 1.5, or 2.5% bases for each diverged monomer template (Additional file 1: Fig. S3b, c). In each case, we can obtain 36 ($4 \times 3 \times 3$) sets of simulated data and we repeated 10 times. Finally, we obtained 360 sets of canonical HOR data and 360 of LN-HOR. The code for generating simulation data is in <https://github.com/xjtu-omics/HiCAT>.

We compared HiCAT with HORmon using default parameter. And the bias measurement is defined as:

$$\text{Bias} = |AN - GT| \quad (5)$$

where AN is the number of annotation HOR number and GT is ground truth HOR number which is 40 HOR units.

Enrichment analysis of HORs in CEN5

We split the CHM13 CEN5 into 25 bins to determine the location specificity of 12, 16, and 20 monomer HOR patterns. Firstly, we randomly generated the same number of 12, 16 and 20 monomer HOR patterns based on uniform distribution across CEN5 with 100 times. Then, we calculated mean (μ) and standard deviation (σ) of the background distribution for each HOR pattern, and calculated the $z = (N_{obv} - \mu)/\sigma$, where N_{obv} is the observed pattern number and p -value is calculated by z based on the standard normal distribution.

Annotation visualization

StainedGlass [22] was used to visualize the TR structures with identity heatmaps, and the window size was set to 2000. We used Gepard [23] to create dot plots. For HiCAT results, within each centromere, we visualized the top five HORs with repeat numbers greater than 10 and reported all detected HORs in the output files.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-02900-5>.

Additional file 1: Supplementary method. Pseudocode for hierarchical tandem repeated mining. **Fig. S1.** Examples of Local-nested HORs (LN-HORs) in CHM13. **Fig. S2.** Two modules in hierarchical tandem repeated mining method. **Fig. S3.** Building simulated datasets. **Fig. S4.** The annotation in human CHM13 CEN19, 9, 10 and 13. **Fig. S5.** Compared HiCAT annotation with previous semi-manual HOR in CEN8. **Fig. S6.** The annotation in human CHM13 CEN5 and 17. **Fig. S7.** Compared with HORmon in human CHM13. **Fig. S8.** Dot plots and similarity heatmap for R18L8 in Arabidopsis thaliana CEN2. **Fig. S9.** Annotation centromere structures in Arabidopsis thaliana CEN2, CEN3, CEN4 and CEN5. **Fig. S10.** The different hypothetical monomer arrangements in HiCAT annotation.

Additional file 2: Table S1. Simulated validation in canonical datasets.

Additional file 3: Table S2. Simulated validation in local nested HOR(LN-HOR) datasets.

Additional file 4: Table S3. Centromere regions in human and Arabidopsis thaliana.

Additional file 5: Table S4. Comparison between HiCAT and previous studies for HORs in human centromeres.

Additional file 6: Table S5. Location bias enrichment test in CEN5.

Additional file 7: Table S6. Comparison annotation coverage and continuity in CHM13 centromeres for HiCAT and HORmon.

Additional file 8: Table S7. HiCAT HORs in Arabidopsis thaliana centromeres.

Additional file 9: Table S8. Runtime for HiCAT in human and Arabidopsis thaliana.

Additional file 10. Review history.

Acknowledgements

The authors thank Yujing Liu and Peng Jia for their important suggestions and feedback.

Review history

The review history is available as Additional file 10.

Peer review information

Andrew Cosgrove was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

Authors' contributions

KY and XY conceived the study. SG, BW, and XZ analyzed the data. SG, XY, and HG developed the program. SG and XY wrote the manuscript. SG completed figures of manuscript. All authors read and approved the final manuscript.

Funding

This work is supported by the National Science Foundation of China (32125009, 62172325, 32070663), the National Key R&D Program of China (2022YFC3400300), the Fundamental Research Funds for the Central Universities, the Natural Science Basic Research Program of Shaanxi (2021GXLH-Z-098), the Science and Technology Project of Xi'an (21RGSF0013), and Key Construction Program of the National "985" Project.

Availability of data and materials

The datasets supported this article including genomes (human CHM13 assembly [24] and *Arabidopsis thaliana* Col-CEN assembly [25]), the centromere coordinates were summarized in Additional file 4: Table S3), HORmon annotation of CHM13 centromeres [26] and T2T Consortium annotation of CHM13 centromeres [27]. Source code of HICAT is available under the GPLv2 license (<https://opensource.org/licenses/GPL-2.0>) from Github [28] and Zenodo [29].

Declarations**Ethics approval and consent to participate**

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 12 July 2022 Accepted: 17 March 2023

Published online: 28 March 2023

References

- McKinley KL, Cheeseman IM. The molecular basis for centromere identity and function. *Nat Rev Mol Cell Biol.* 2016;17:16–29.
- Henikoff S, Ahmad K, Malik HS. The centromere paradox: stable inheritance with rapidly evolving DNA. *Science.* 2001;293:1098–102.
- McNulty SM, Sullivan BA. Alpha satellite DNA biology: finding function in the recesses of the genome. *Chromosome Res.* 2018;26:115–38.
- Dvorkina T, Bzikadze AV, Pevzner PA. The string decomposition problem and its applications to centromere analysis and assembly. *Bioinformatics.* 2020;36:i93–101.
- Bzikadze AV, Pevzner PA. Automated assembly of centromeres from ultra-long error-prone reads. *Nat Biotechnol.* 2020;38:1309–16.
- Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, Ebler J, Functammasan A, Kolesnikov A, Olson ND, et al. Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat Biotechnol.* 2019;37:1155–62.
- Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze AV, Mikheenko A, Vollger MR, Altemose N, Uralsky L, Gershman A, et al. The complete sequence of a human genome. *Science.* 2022;376:44–53.
- Naish M, Alonge M, Wlodzimierz P, Tock AJ, Abramson BW, Schmucker A, Mandakova T, Jamge B, Lambing C, Kuo P, et al. The genetic and epigenetic landscape of the *Arabidopsis* centromeres. *Science.* 2021;374:eabi7489.
- Song JM, Xie WZ, Wang S, Guo YX, Koo DH, Kudrna D, Gong C, Huang Y, Feng JW, Zhang W, et al. Two gap-free reference genomes and a global view of the centromere architecture in rice. *Mol Plant.* 2021;14:1757–67.
- Dvorkina T, Kunyavskaya O, Bzikadze AV, Alexandrov I, Pevzner PA. CentromereArchitect: inference and analysis of the architecture of centromeres. *Bioinformatics.* 2021;37:i196–204.
- Altemose N, Logsdon GA, Bzikadze AV, Sidhwani P, Langley SA, Caldas GV, Hoyt SJ, Uralsky L, Ryabov FD, Shew CJ, et al. Complete genomic and epigenetic maps of human centromeres. *Science.* 2022;376:eab14178.
- Shepelev VA, Uralsky LI, Alexandrov AA, Yurov YB, Rogaev EI, Alexandrov IA. Annotation of suprachromosomal families reveals uncommon types of alpha satellite organization in pericentromeric regions of hg38 human genome assembly. *Genom Data.* 2015;5:139–46.
- Uralsky LI, Shepelev VA, Alexandrov AA, Yurov YB, Rogaev EI, Alexandrov IA. Classification and monomer-by-monomer annotation dataset of suprachromosomal family 1 alpha satellite higher-order repeats in hg38 human genome assembly. *Data Brief.* 2019;24:103708.
- Kunyavskaya O, Dvorkina T, Bzikadze AV, Alexandrov IA, Pevzner PA. Automated annotation of human centromeres with HORmon. *Genome Res.* 2022;32:1137–51.

15. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech*. 2008;2008:P10008.
16. Traag VA, Waltman L, van Eck NJ. From Louvain to Leiden: guaranteeing well-connected communities. *Sci Rep*. 2019;9:5233.
17. Talbert PB, Henikoff S. What makes a centromere? *Exp Cell Res*. 2020;389:111895.
18. Logsdon GA, Vollger MR, Hsieh P, Mao Y, Liskovych MA, Koren S, Nurk S, Mercuri L, Dishuck PC, Rhie A, et al. The structure, function and evolution of a complete human chromosome 8. *Nature*. 2021;593:101–7.
19. Alexandrov I, Kazakov A, Tumeneva I, Shepelev V, Yurov Y. Alpha-satellite DNA of primates: old and new families. *Chromosoma*. 2001;110:253–66.
20. Sevim V, Bashir A, Chin CS, Miga KH. Alpha-CENTAURI: assessing novel centromeric repeat sequence variation with long read sequencing. *Bioinformatics*. 2016;32:1921–4.
21. Suzuki Y, Myers EW, Morishita S. Rapid and ongoing evolution of repetitive sequence structures in human centromeres. *Sci Adv*. 2020;6:eabd9230.
22. Vollger MR, Kerpedjiev P, Phillippy AM, Eichler EE. StainedGlass: Interactive visualization of massive tandem repeat structures with identity heatmaps. *Bioinformatics*. 2022;38:2049–51.
23. Krumsiek J, Arnold R, Rattei T. Gepard: a rapid and sensitive tool for creating dotplots on genome scale. *Bioinformatics*. 2007;23:1026–8.
24. Nurk S, Koren S, Rhie A, Rautiainen M, Bizikadze AV, Mikheenko A, Vollger MR, Altemose N, Uralsky L, Gershman A, et al. The complete sequence of a human genome. *Datasets*. Github. <https://github.com/marbl/CHM13> (2022).
25. Naish M, Alonge M, Wlodzimierz P, Tock AJ, Abramson BW, Schmucker A, Mandakova T, Jamge B, Lambing C, Kuo P, et al. The genetic and epigenetic landscape of the Arabidopsis centromeres. *Datasets*. Github. <https://github.com/schatzlab/Col-CEN> (2021).
26. Kunyavskaya O, Dvorkina T, Bizikadze AV, Alexandrov IA, Pevzner PA. Automated annotation of human centromeres with HORmon. *Datasets*. Figshare. <https://figshare.com/articles/dataset/HORmon/16755097/2> (2022).
27. Altemose N, Logsdon GA, Bizikadze AV, Sidhwani P, Langley SA, Caldas GV, Hoyt SJ, Uralsky L, Ryabov FD, Shew CJ, et al. Complete genomic and epigenetic maps of human centromeres. *Datasets*. UCSC browser. http://t2t.gi.ucsc.edu/chm13/hub/t2t-chm13-v1.0/alphaSatHOR/ASat_HOR.bigBed (2022).
28. Gao S, Yang X, Guo H, Zhao X, Wang B, Ye K. HiCAT: A tool for automatic annotation of centromere structure. *Github*. <https://github.com/xjtu-omics/HiCAT> (2022).
29. Gao S, Yang X, Guo H, Zhao X, Wang B, Ye K. HiCAT: a tool for automatic annotation of centromere structure. *Zenodo*. <https://doi.org/10.5281/zenodo.7260510> (2022).

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

