



OPEN

Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes

Kishwar Shafin^{1,11}, Trevor Pesout^{1,11}, Ryan Lorig-Roach^{1,11}, Marina Haukness^{1,11}, Hugh E. Olsen^{1,11}, Colleen Bosworth¹, Joel Armstrong¹, Kristof Tigyi^{1,2}, Nicholas Maurer¹, Sergey Koren³, Fritz J. Sedlazeck⁴, Tobias Marschall⁵, Simon Mayes⁶, Vania Costa⁶, Justin M. Zook⁷, Kelvin J. Liu⁸, Duncan Kilburn⁸, Melanie Sorensen⁹, Katy M. Munson⁹, Mitchell R. Vollger⁹, Jean Monlong¹, Erik Garrison¹, Evan E. Eichler^{2,9}, Sofie Salama^{1,2}, David Haussler^{1,2}, Richard E. Green¹, Mark Akeson¹, Adam Phillippy³, Karen H. Miga¹, Paolo Carnevali¹⁰✉, Miten Jain¹✉ and Benedict Paten¹✉

De novo assembly of a human genome using nanopore long-read sequences has been reported, but it used more than 150,000 CPU hours and weeks of wall-clock time. To enable rapid human genome assembly, we present Shasta, a de novo long-read assembler, and polishing algorithms named MarginPolish and HELEN. Using a single PromethION nanopore sequencer and our toolkit, we assembled 11 highly contiguous human genomes de novo in 9 d. We achieved roughly 63× coverage, 42-kb read N50 values and 6.5× coverage in reads >100 kb using three flow cells per sample. Shasta produced a complete haploid human genome assembly in under 6 h on a single commercial compute node. MarginPolish and HELEN polished haploid assemblies to more than 99.9% identity (Phred quality score QV = 30) with nanopore reads alone. Addition of proximity-ligation sequencing enabled near chromosome-level scaffolds for all 11 genomes. We compare our assembly performance to existing methods for diploid, haploid and trio-binned human samples and report superior accuracy and speed.

Reference-based methods such as GATK¹ can infer human variations from short-read sequences, but the results only cover ~90% of the reference human genome assembly^{2,3}. These methods are accurate with respect to single-nucleotide variants and short insertions and deletions (indels) in this mappable portion of the reference genome⁴. However, it is difficult to use short reads for de novo genome assembly⁵, to discover structural variations (SVs)^{6,7} (including large indels and base-level resolved copy number variations), or to resolve phasing relationships without exploiting transmission information or haplotype panels⁸.

Third generation sequencing technologies, including linked-reads^{9–11} and long-read technologies^{12,13}, overcome the fundamental limitations of short-read sequencing for genome inference. In addition to increasingly being used in reference guided methods^{2,14–16}, long-read sequences can generate highly contiguous de novo genome assemblies¹⁷.

Nanopore sequencing, as commercialized by Oxford Nanopore Technologies (ONT), is particularly useful for de novo genome assembly because it can produce high yields of very long 100+ kilobase (kb) reads¹⁸. Very long reads hold the promise of facilitating contiguous, unbroken assembly of the most challenging regions of the human genome, including centromeric satellites, acrocentric

short arms, ribosomal DNA arrays and recent segmental duplications^{19–21}. The de novo assembly of a nanopore sequencing based human genome has been reported¹⁸. This earlier effort needed 53 ONT MinION flow cells and the assembly required more than 150,000 CPU hours and weeks of wall-clock time, quantities that are unfeasible for high throughput human genome sequencing efforts.

To enable easy, cheap and fast de novo assembly of human genomes we developed a toolkit for nanopore data assembly and polishing that is orders of magnitude faster than state-of-the-art methods. We use a combination of nanopore and proximity-ligation (HiC) sequencing⁹ and our toolkit, and we report improvements in human genome sequencing coupled with reduced time, labor and cost.

Results

Eleven human genomes sequenced in 9 d. We selected for sequencing 11, low-passage (six passages), human cell lines of the offspring of parent-child trios from the 1,000 Genomes Project²² and genome-in-a-bottle (GIAB)²³ sample collections. Samples were selected to maximize captured allelic diversity (see Methods).

We carried out PromethION nanopore sequencing and HiC Illumina sequencing for the 11 genomes. We used three flow cells per genome, with each flow cell receiving a nuclease flush every

¹UC Santa Cruz Genomics Institute, Santa Cruz, CA, USA. ²Howard Hughes Medical Institute, University of California, Santa Cruz, CA, USA. ³Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, Bethesda, MD, USA. ⁴Baylor College of Medicine, Human Genome Sequencing Center, Houston, TX, USA. ⁵Max Planck Institute for Informatics, Saarbrücken, Germany. ⁶Oxford Nanopore Technologies, Oxford, UK. ⁷National Institute of Standards and Technology, Gaithersburg, MD, USA. ⁸Circulomics Inc., Baltimore, MD, USA. ⁹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA. ¹⁰Chan Zuckerberg Initiative, Redwood City, CA, USA. ¹¹These authors contributed equally: Kishwar Shafin, Trevor Pesout, Ryan Lorig-Roach, Marina Haukness, Hugh E. Olsen. ✉e-mail: paolo@chanzuckerberg.com; miten@soe.ucsc.edu; bpaten@ucsc.edu

20–24 h. This flush removed long DNA fragments that could cause the pores to become blocked over time. Each flow cell received a fresh library of the same sample after the nuclease flush. A total of two nuclease flushes were performed per flow cell, and each flow cell received a total of three sequencing libraries. We used Guppy v.2.3.5 with the high accuracy flipflop model for basecalling (see Methods).

Nanopore sequencing of all 11 genomes took 9 d and produced 2.3 terabases (Tb) of sequence. We ran up to 15 flow cells in parallel during these sequencing runs. Results are shown in Fig. 1 and Supplementary Tables 1–3. Nanopore sequencing yielded an average of 69 gigabases (Gb) per flow cell, with the total throughput per individual genome ranging between 48× (158 Gb) and 85× (280 Gb) coverage per genome (Fig. 1a). The read N50 values for the sequencing runs ranged between 28 and 51 kb (Fig. 1b). (An N50 value is a weighted median; it is the length of the sequence in a set for which all sequences of that length or greater sum to 50% of the set's total size.) We aligned nanopore reads to the human reference genome (GRCh38) and calculated their alignment identity to assess sequence quality (see Methods). We observed that the median and modal alignment identity was 90 and 93%, respectively (Fig. 1c). The sequencing data per individual genome included an average of 55× coverage arising from >10-kb reads and 6.5× coverage from >100-kb reads (Fig. 1d). This was in large part due to size selection that yielded an enrichment of reads longer than 10 kb. To test the generality of our sequencing methodology for other samples, we sequenced high-molecular weight DNA isolated from a human saliva sample using identical sample preparation. The library was run on a MinION (roughly one-sixth the throughput of a ProMethION flow cell) and yielded 11 Gb of data at a read N50 of 28 kb (Supplementary Table 4), extrapolating both are within the lower range achieved with cell-line derived DNA.

Shasta assembler for long sequence reads. Shasta was designed to be orders of magnitude faster and cheaper at assembling a human-scale genome from nanopore reads than the Canu assembler used in our earlier work¹⁸. During most Shasta assembly phases, reads are stored in a homopolymer-compressed form using run-length encoding (RLE)^{24–26}. In this form, identical consecutive bases are collapsed, and the base and repeat count are stored. For example, GATTTACCA would be represented as (GATACA, 113121). This representation is insensitive to errors in the length of homopolymer runs, thereby addressing the dominant error mode for Oxford Nanopore reads¹². As a result, assembly noise due to read errors is decreased, and notably higher identity alignments are facilitated (Fig. 1e). A marker representation of reads is also used, in which each read is represented as the sequence of occurrences of a predetermined, fixed subset of short *k*-mers (marker representation) in its run-length representation. A modified MinHash^{27,28} scheme is used to find candidate pairs of overlapping reads, using as MinHash features consecutive occurrences of *m* markers (default *m* = 4). Optimal alignments in marker representation are computed for all candidate pairs. The computation of alignments in marker representation is very efficient, particularly as various banded heuristics are used. A marker graph is created in which each vertex represents a marker found to be aligned in a set of several reads. The marker graph is used to assemble sequence after undergoing a series of simplification steps. The assembler runs on a single machine with a large amount of memory (typically 1–2 Tb for a human assembly). All data structures are kept in memory, and no disk I/O takes place except for initial loading of the reads and final output of assembly results.

Benchmarking Shasta. We compared Shasta to three contemporary assemblers: Wtdbg2 (ref. 29), Flye³⁰ and Canu³¹. We ran all four assemblers on available read data from two diploid human samples,

HG00733 and HG002, and one haploid human sample, CHM13. HG00733 and HG002 were part of our collection of 11 samples, and data for CHM13 came from the T2T consortium³².

Canu consistently produced the most contiguous assemblies, with contig NG50 values of 40.6, 32.3 and 79.5 Mb, for samples HG00733, HG002 and CHM13, respectively (Fig. 2a). (NG50 is similar to N50, but for 50% of the estimated genome size.) Flye was the second most contiguous, with contig NG50 values of 25.2, 25.9 and 35.3 Mb, for the same samples. Shasta was next with contig NG50 values of 21.1, 20.2 and 41.1 Mb. Wtdbg2 produced the least contiguous assemblies, with contig NG50 values of 15.3, 13.7 and 14.0 Mb.

Conversely, aligning the samples to GRCh38 and evaluating with QUAST³³, Shasta had between 4.2 and 6.5× fewer disagreements (locations where the assembly contains a breakpoint with respect to the reference assembly) per assembly than the other assemblers (Supplementary Table 5). Breaking the assemblies at these disagreements and unaligned regions with respect to GRCh38, we observe much smaller absolute variation in contiguity (Fig. 2b and Supplementary Table 5). However, a substantial fraction of the identified disagreements likely reflect true SVs with respect to GRCh38. To address this, we discounted disagreements within chromosome Y, centromeres, acrocentric chromosome arms, QH-regions and known recent segmental duplications (all of which are enriched in SVs^{34,35}); in the case of HG002, we further excluded a set of known SVs³⁶. We still observe between 1.2× and 2× fewer disagreements in Shasta relative to Canu and Wtdbg2, and comparable results against Flye (Fig. 2c and Supplementary Table 6). To account for differences in the fraction of the genomes assembled, we analyzed disagreements contained within the intersection of all the assemblies (that is, in regions where all assemblers produced a unique assembled sequence). This produced results highly consistent with the previous analysis and suggests Shasta and Flye have the lowest and comparable rates of misassembly (Methods, see Supplementary Table 7). Finally, we used QUAST to calculate disagreements between the T2T Consortium's chromosome X assembly, a highly curated, validated assembly³² and the subset of each CHM13 assembly mapping to it; Shasta has two to 17 times fewer disagreements than the other assemblers while assembling almost the same fraction of the assembly (Supplementary Table 8).

Canu consistently assembled the largest genomes (average 2.91 Gb), followed by Flye (average 2.83 Gb), Wtdbg2 (average 2.81 Gb) and Shasta (average 2.80 Gb). Due to their similarity, we would expect the most of these assembled sequences to map to another human genome. Discounting unmapped sequence, the differences are smaller: Canu produced an average of 2.86 Gb of mapped sequence per assembly, followed by Shasta (average 2.79 Gb), Flye (average 2.78 Gb) and Wtdbg2 (average 2.76 Gb) (Fig. 2d, see Methods). This analysis supports the notion that Shasta is currently relatively conservative versus its peers, producing the highest ratio of directly mapped assembly per sample.

For HG00733 and CHM13 we examined a library of bacterial artificial chromosome (BAC) assemblies (Methods). The BACs were largely targeted at known segmental duplications (473 of 520 BACs lie within 10 kb of a known duplication). Examining the subset of BACs for CHM13 and HG00733 that map to unique regions of GRCh38 (see Methods), we find Shasta contiguously assembles all 47 BACs, with Flye performing similarly (Supplementary Table 9). In the full set, we observe that Canu (411) and Flye (282) contiguously assemble a larger subset of the BACs than Shasta (132) and Wtdbg2 (108), confirming the notion that Shasta is relatively conservative in these duplicated regions (Supplementary Table 10). Examining the fraction of contiguously assembled BACs of all BACs represented in each assembly we can measure an aspect of assembly correctness. In this regard Shasta (97%) produces a much higher percentage of correct BACs in duplicated regions versus its peers

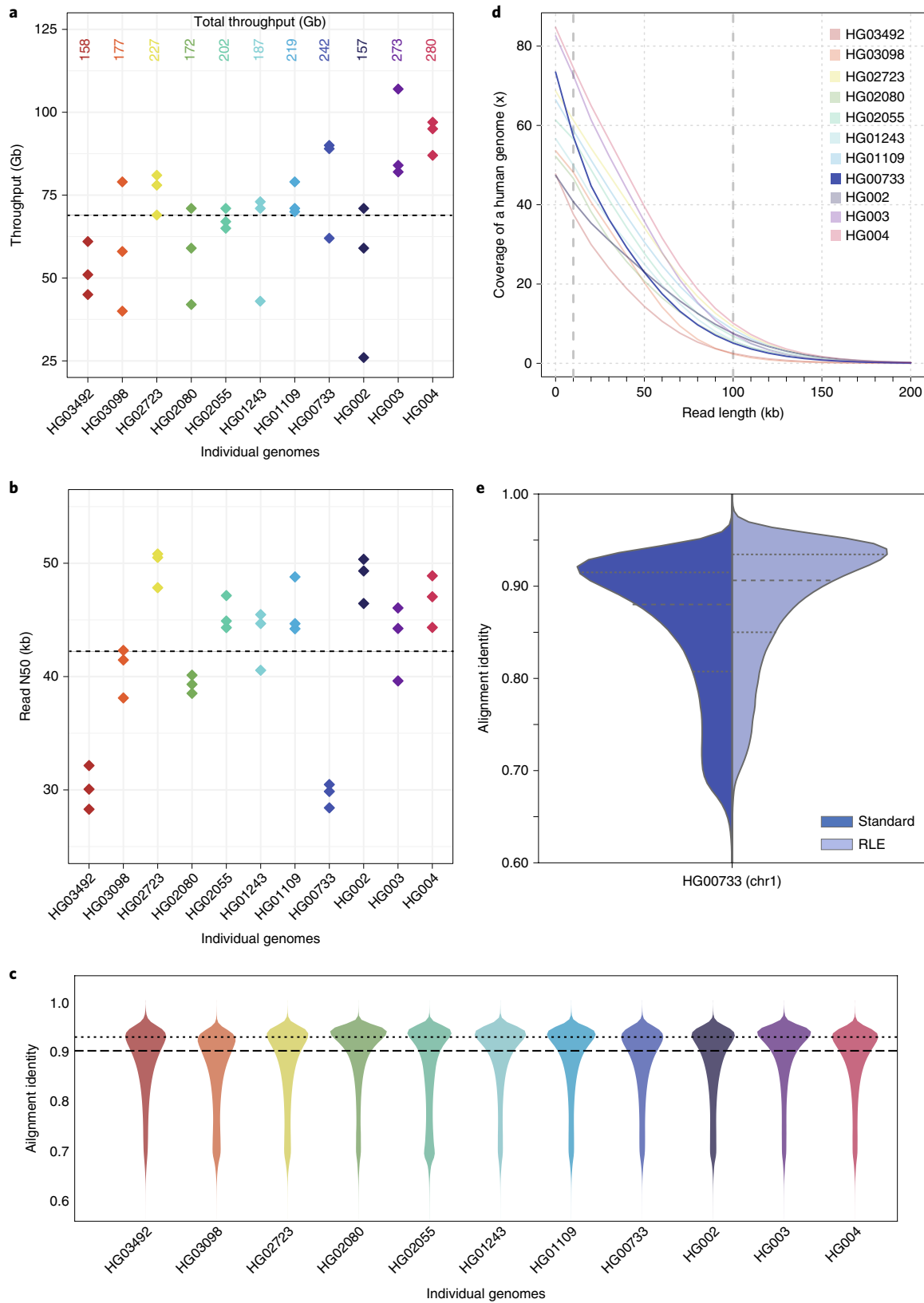


Fig. 1 | Nanopore sequencing data. **a**, Throughput in gigabases from each of three flow cells for 11 samples, with total throughput at top. Each point is a flow cell. **b**, Read N50 values for each flow cell. Each point is a flow cell. **c**, Alignment identities against GRCh38. Medians in **a–c** shown by dashed lines, dotted line in **c** is the mode. Each line is a single sample comprising three flow cells. **d**, Genome coverage as a function of read length. Dashed lines indicate coverage at 10 and 100 kb. HG00733 is accentuated in dark blue as an example. Each line is a single sample comprising three flow cells. **e**, Alignment identity for standard and RLE reads. Data for HG00733 chromosome 1 flow cell 1 are shown (4.6 Gb raw sequence). Dashed lines denote quartiles.

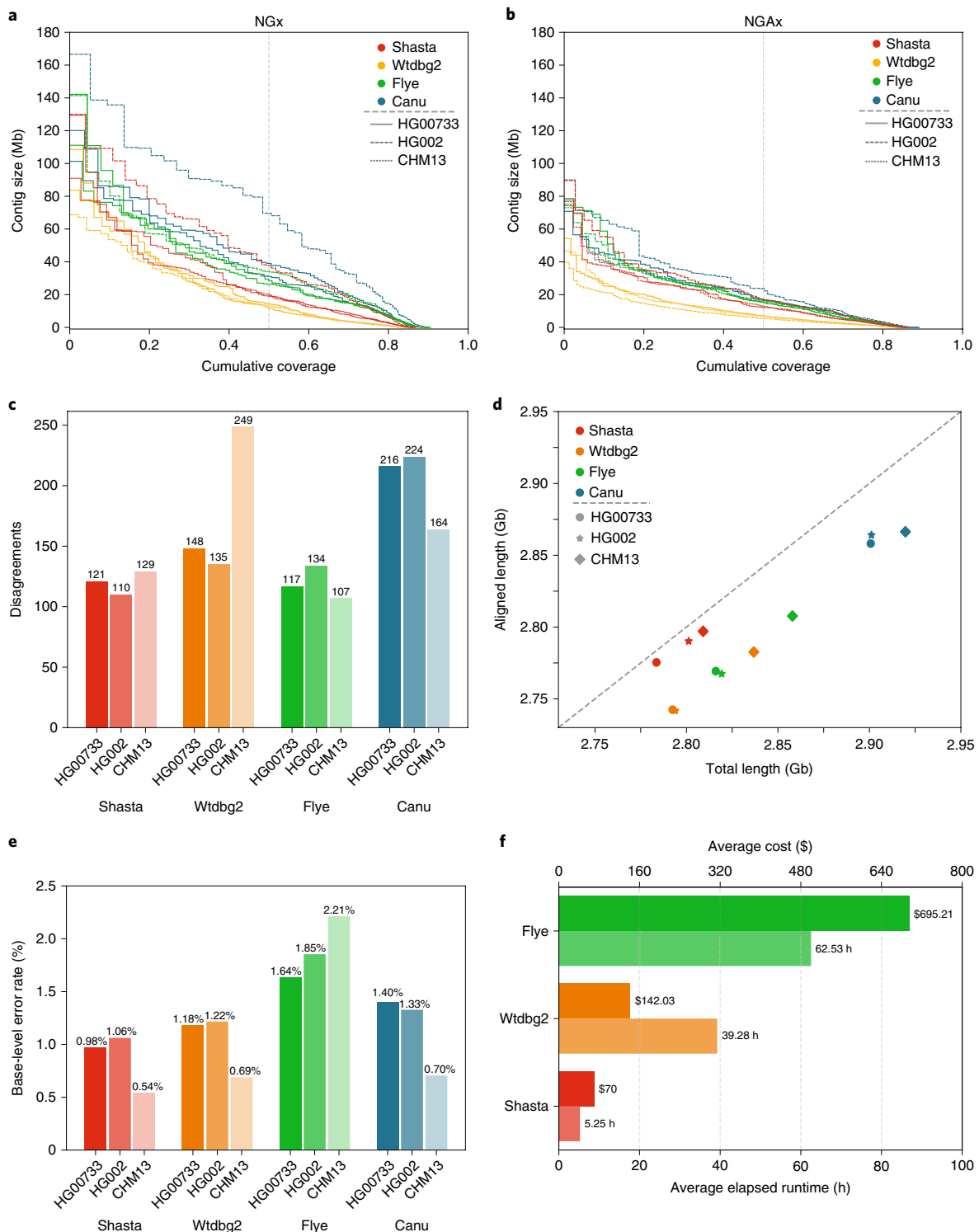


Fig. 2 | Assembly metrics for Shasta, Wtdbg2, Flye and Canu before polishing. **a**, NGx plot showing contig length distribution. The intersection of each line with the dashed line is the NG50 for that assembly. **b**, NGAx plot showing the distribution of aligned contig lengths. Each horizontal line represents an aligned segment of the assembly unbroken by a disagreement or unmappable sequence with respect to GRCh38. The intersection of each line with the dashed line is the aligned NGA50 for that assembly. **c**, Assembly disagreement counts for regions outside centromeres, segmental duplications and, for HG002, known SVs. **d**, Total generated sequence length versus total aligned sequence length (against GRCh38). **e**, Balanced base-level error rates for assembled sequences. **f**, Average runtime and cost for assemblers (Canu not shown).

(Canu 92%, Flye 87%, Wtdbg2 88%). In the intersected set of BACs attempted by all assemblers (Supplementary Table 11), Shasta, 100%; Flye, 100%; Canu, 98.50% and Wtdbg2, 90.80% all produce comparable results.

Shasta produced the most base-level accurate assemblies (average balanced error rate 0.98% on diploid and 0.54% on haploid), followed by Wtdbg2 (1.18% on diploid and 0.69% on haploid), Canu (1.40% on diploid and 0.71% on haploid) and Flye (1.64% on diploid

and 2.21% on haploid) (Fig. 2e, see Methods and Supplementary Table 12). We also calculated the base-level accuracy in regions covered by all the assemblies and observe results consistent with the whole-genome assessment (Supplementary Table 13).

Shasta, Wtdbg2 and Flye were run on a commercial cloud, allowing us to reasonably compare their cost and runtime (Fig. 2e, see Methods). Shasta took an average of 5.25 h to complete each assembly at an average cost of US\$70 per sample. In contrast, Wtdbg2 took 7.5× longer and cost 3.7× as much, and Flye took 11.9× longer and cost 9.9× as much. Due to the anticipated cost and complexity of porting it to Amazon Web Services (AWS), the Canu assemblies were run on a large, institutional compute cluster, consuming up to US\$19,000 (estimated) of compute and took around 4–5 d per assembly (Methods, see Supplementary Tables 14 and 15).

To assess the use of using Shasta for SV characterization we created a workflow to extract putative heterozygous SVs from Shasta assembly graphs (Methods). Extracting SVs from an assembly graph for HG002, the length distribution of indels shows the characteristic spikes for known retrotransposon lengths (Supplementary Fig. 1). Comparing these SVs to the high-confidence GIAB SV set we find good concordance, with a combined F1 score of 0.68 (Supplementary Table 16).

Contiguously assembling major histocompatibility complex (MHC) haplotypes. The MHC region is difficult to resolve using short reads due to its repetitive and highly polymorphic nature³⁷, and recent efforts to apply long-read sequencing to this problem have shown promise^{18,38}. We analyzed the assemblies of CHM13 and HG00733 to see if they spanned the MHC region. For the haploid assemblies of CHM13 we find MHC is entirely spanned by a single contig in all four assemblers' output, and most closely resembles the GL000251.2 haplogroup among those provided in GRCh38 (Fig. 3a, Supplementary Fig. 2 and Supplementary Table 17). In the diploid assembly of HG00733 two contigs span most of the MHC for Shasta and Flye, while Canu and Wtdbg2 span the region with one contig (Fig. 3b and Supplementary Fig. 3). However, we note that all these chimeric diploid assemblies lead to sequences that do not closely resemble any haplogroup (Methods).

To attempt to resolve haplotypes of HG00733 we used trio-binning³⁹ to partition the reads for HG00733 into two sets based on likely maternal or paternal lineage and assembled the haplotypes (Methods). For all assemblers and each haplotype assembly, the global contiguity worsened substantially (as the available read data coverage was approximately halved and, further, not all reads could be partitioned), but the resulting disagreement count decreased (Supplementary Table 18). When using haploid trio-binned assemblies, the MHC was spanned by a single contig for the maternal haplotype (Fig. 3c, Supplementary Fig. 4 and Supplementary Table 19), with high identity to GRCh38 and having the greatest contiguity and identity with the GL000255.1 haplotype. For the paternal haplotype, low coverage led to discontinuities (Fig. 3d) breaking the region into three contigs.

Deep neural network-based polishing for long-read assemblies.

We developed a deep neural network-based consensus sequence polishing pipeline designed to improve the base-level quality of the initial assembly. The pipeline consists of two modules: MarginPolish and the homopolymer encoded long-read error-corrector for Nanopore (HELEN). MarginPolish uses a banded form of the forward-backward algorithm on a pairwise hidden Markov model (pair-HMM) to generate pairwise alignment statistics from the RLE alignment of each read to the assembly. From these statistics, MarginPolish generates a weighted RLE partial order alignment (POA)⁴⁰ graph that represents potential alternative local assemblies. MarginPolish iteratively refines the assembly using this RLE POA, and then outputs the final summary graph for consumption

by HELEN. HELEN uses a multi-task recurrent neural network (RNN)⁴¹ that takes the weights of the MarginPolish RLE POA graph to predict a nucleotide base and run length for each genomic position. The RNN takes advantage of contextual genomic features and associative coupling of the POA weights to the correct base and run length to produce a consensus sequence with higher accuracy.

To demonstrate the effectiveness of MarginPolish and HELEN, we compared them with the state-of-the-art nanopore assembly polishing workflow: four iterations of Racon polishing⁴² followed by Medaka. MarginPolish is analogous in function to Racon, both using pair-HMM-based methods for alignment and POA graphs for initial refinement. Similarly, HELEN is analogous to Medaka, in that both use a deep neural network and both work from summary statistics of reads aligned to the assembly.

Figure 4a and Supplementary Tables 20–22 detail error rates for the four methods performed on the HG00733 and CHM13 Shasta assemblies (see Methods) using Pomoxis. For the diploid HG00733 sample MarginPolish and HELEN achieve a balanced error rate of 0.388% (Phred quality score QV = 24.12), compared to 0.455% (QV = 23.42) by Racon and Medaka. For both polishing pipelines, a notable fraction of these errors are likely due to true heterozygous variations. For the haploid CHM13 we restrict comparison to the highly curated X chromosome sequence provided by the T2T consortium³². We achieve a balanced error rate of 0.064% (QV = 31.92), compared to Racon and Medaka's 0.110% (QV = 29.59).

For all assemblies, errors were dominated by indel errors; for example, substitution errors are 3.16 and 2.9 times fewer than indels in the polished HG00733 and CHM13 assemblies, respectively. Many of these errors relate to homopolymer length confusion; Fig. 4b analyzes the homopolymer error rates for various steps of the polishing workflow for HG00733. Each panel shows a heatmap with the true length of the homopolymer run on the y axis and the predicted run length on the x axis, with the color describing the likelihood of predicting each run length given the true length. Note that the dispersion of the diagonal steadily decreases. The vertical streaks at high run lengths in the MarginPolish and HELEN confusion matrix are the result of infrequent numerical and encoding artifacts (see Methods and Supplementary Fig. 5).

Figure 4c and Supplementary Table 23 show the overall error rate after running MarginPolish and HELEN on HG00733 assemblies generated by different assembly tools, demonstrating that they can be usefully employed to polish assemblies generated by other tools.

To investigate the benefit of using short reads for further polishing, we polished chromosome X of the CHM13 Shasta assembly after MarginPolish and HELEN using 10X Chromium reads with the Pilon polisher⁴³. This led to a roughly twofold reduction in base errors, increasing the Phred quality score from roughly QV = 32 (after polishing with MarginPolish and HELEN) to around QV = 36 (Supplementary Table 24). Notably, attempting to use Pilon polishing on the raw Shasta assembly resulted in much poorer results (QV = 24).

Figure 4d and Supplementary Table 25 describe average run-times and costs for the methods (see Methods). MarginPolish and HELEN cost a combined US\$107 and took 29 h of wall-clock time on average, per sample. In comparison Racon and Medaka cost US\$621 and took 142 wall-clock hours on average, per sample. To assess single-region performance we additionally ran the two polishing workflows on a single contig (roughly 1% of the assembly size), MarginPolish/HELEN was three times faster than Racon (1×)/Medaka (Supplementary Table 26).

Long-read assemblies contain nearly all human coding genes.

To evaluate the accuracy and completeness of an assembled transcriptome we ran the Comparative Annotation Toolkit⁴⁴, which can annotate a genome assembly using the human GENCODE⁴⁵

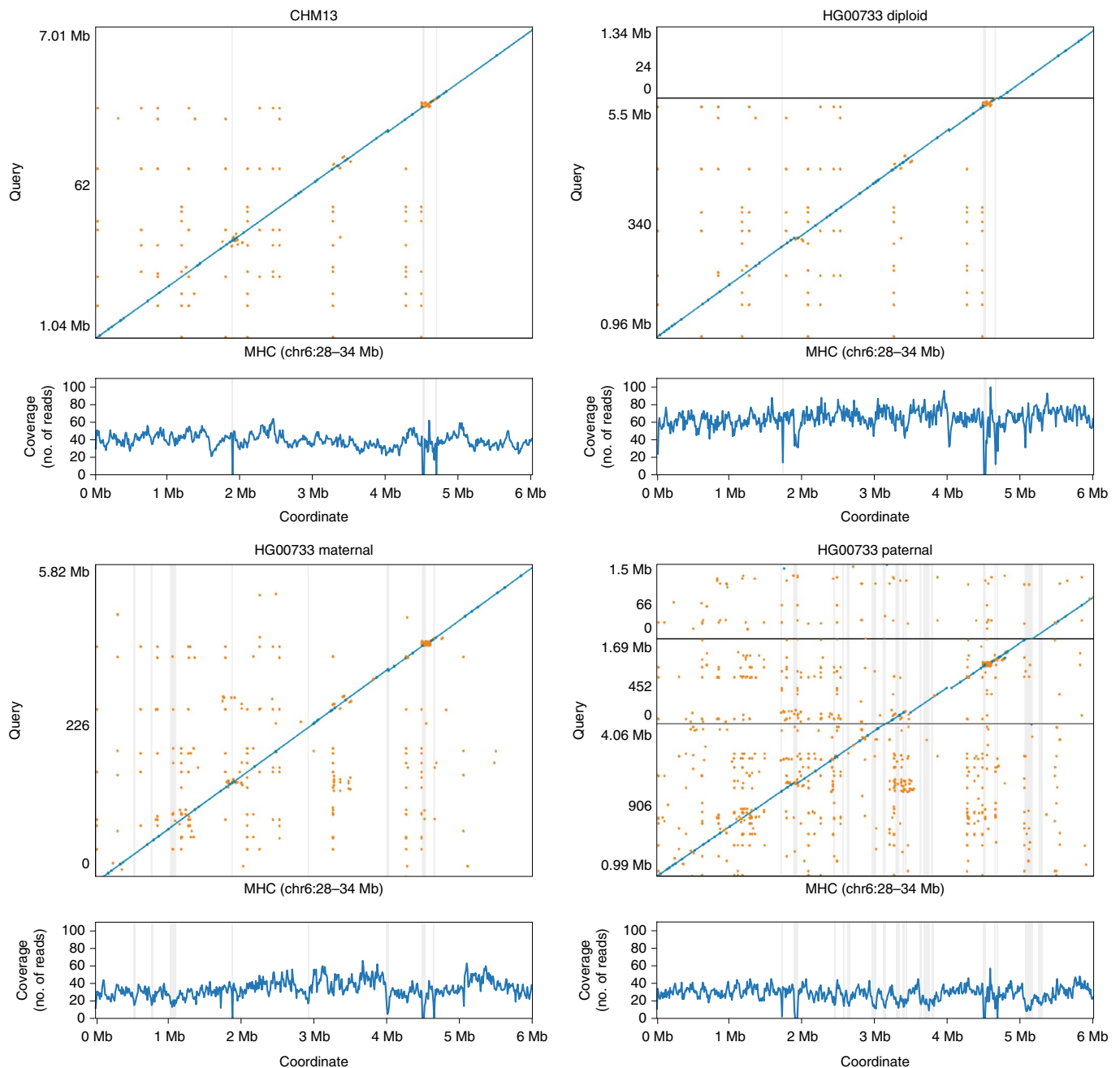


Fig. 3 | Shasta MHC assemblies compared with the reference human genome. Unpolished Shasta assembly for CHM13 and HG00733, including HG00733 trio-binned maternal and paternal assemblies. Shaded gray areas are regions in which coverage (as aligned to GRCh38) drops below 20. Horizontal black lines indicate contig breaks. Blue and green describe unique alignments (aligning forward and reverse, respectively) and orange describes multiple alignments.

reference human gene set (Table 1, Methods and Supplementary Tables 27–30).

For the HG00733 and CHM13 samples we found that Shasta assemblies polished with MarginPolish and HELEN contained nearly all human protein coding genes, having, respectively, an identified ortholog for 99.23% (152 missing) and 99.11% (175 missing) of these genes. Using the restrictive definition that a coding gene is complete in the assembly only if it is assembled across its full length, contains no frameshifts and retains the original intron–exon structure, we found that 68.07% and 74.20% of genes, respectively, were complete in the HG00733 and CHM13 assemblies. Polishing the Shasta assemblies alternatively with the

Racon–Medaka pipeline achieved similar but uniformly less complete results.

Comparing the MarginPolish and HELEN polished assemblies for HG00733 generated with Flye, Canu and Wtdbg2 to the similarly polished Shasta assembly we found that Canu had the fewest missing genes (just 51), but that Flye, followed by Shasta, had the most complete genes. Wtdbg2 was clearly an outlier, with notably larger numbers of missing genes (506). For comparison we additionally ran BUSCO⁴⁶ using the eukaryote set of orthologs on each assembly, a smaller set of 303 expected single-copy genes (Supplementary Tables 31 and 32). We find comparable performance between the assemblies, with small

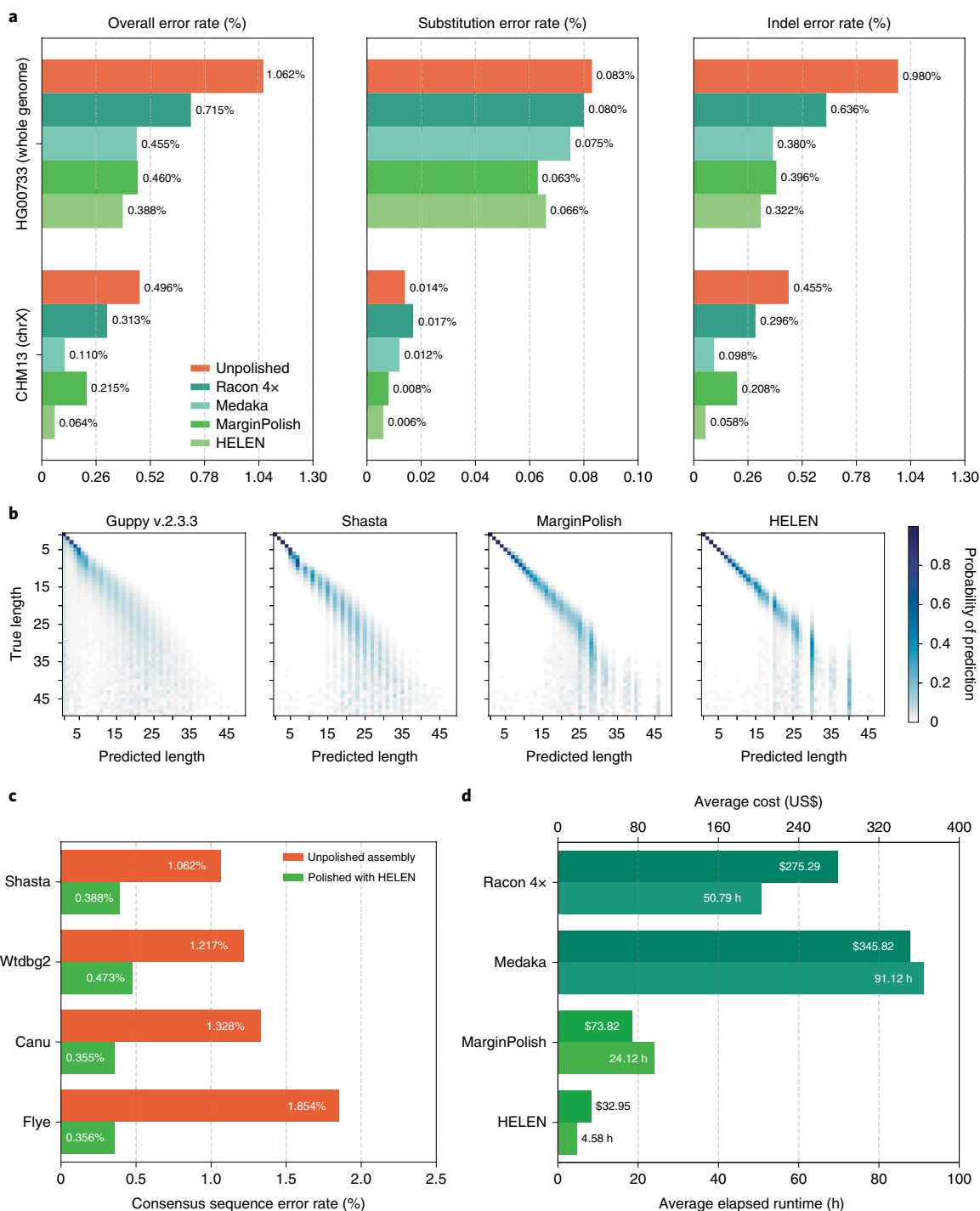


Fig. 4 | Polishing assembled genomes. a, Balanced error rates for the four methods on HG00733 and CHM13. **b**, Row-normalized heatmaps describing the predicted run lengths (x axis) given true run lengths (y axis) for four steps of the pipeline on HG00733. Guppy v.2.3.3 was generated from 3.7 Gb of RLE sequence. Shasta, MarginPolish and HELEN were generated from whole assemblies aligned to their respective truth sequences. **c**, Error rates for MarginPolish and HELEN on four assemblies. **d**, Average runtime and cost.

differences largely recapitulating the pattern observed by the larger CAT analysis.

Comparison of Shasta and PacBio HiFi assemblies. We compared the CHM13 Shasta assembly polished using MarginPolish and HELEN with the recently released Canu assembly of CHM13 using PacBio HiFi reads⁴⁷. HiFi reads are based on circular consensus sequencing technology that delivers substantially lower error rates. The HiFi assembly has a lower NG50 (29.0 versus 41.0 megabase

(Mb)) than the Shasta assembly (Supplementary Fig. 6). Consistent with our other comparisons to Canu, the Shasta assembly also contains a much lower disagreement count relative to GRCh38 (1073) than the Canu-based HiFi assembly (8,469), a difference that remains after looking only at disagreements within the intersection of the assemblies (380 versus 594). The assemblies have an almost equal NGax (~20.0 Mb), but the Shasta assembly covers a smaller fraction of GRCh38 (95.28 versus 97.03%) (Supplementary Fig. 7 and Supplementary Table 33). Predictably, the HiFi assembly

Table 1 | CAT transcriptome analysis of human protein coding genes for HG00733 and CHM13

Sample	Assembler	Polisher	Genes found (%)	Missing genes	Complete genes (%)
HG00733	Canu	HELEN	99.741	51	67.038
	Flye	HELEN	99.405	117	71.768
	Wtdbg2	HELEN	97.429	506	66.143
	Shasta	HELEN	99.228	152	68.069
	Shasta	Medaka	99.141	169	66.27
CHM13	Shasta	HELEN	99.111	175	74.202
	Shasta	Medaka	99.035	190	73.836

has a higher Phred quality score than the polished Shasta assembly (QV = 41 versus QV = 32).

Scaffolding to near chromosome scale. To achieve chromosome length sequences, we scaffolded all of the polished Shasta assemblies with HiC proximity-ligation data using HiRise⁴⁸ (see Methods and Fig. 5a). On average, 891 joins were made per assembly. This increased the scaffold NG50 values to near chromosome scale, with a median of 129.96 Mb, as shown in Fig. 5a, with additional assembly metrics in Supplementary Table 36. Proximity-ligation data can also be used to detect misjoins in assemblies. In all 11 Shasta assemblies, no breaks to existing contigs were made while running HiRise to detect potential misjoins. Aligning HG00733 to GRCh38, we find no notable rearrangements and all chromosomes are spanned by one or a few contigs (Fig. 5b), with the exception of chrY, which is absent because HG00733 is female. Similar results were observed for HG002 (Supplementary Fig. 8).

Discussion

With sequencing efficiency for long reads improving, computational considerations are paramount in determining overall time, cost and quality. Simply put, large genome de novo assembly will not become ubiquitous if the requirements are weeks of assembly time on large computational clusters. We present three new methods that provide a pipeline for the rapid assembly of long nanopore reads. Shasta can produce a draft human assembly in around 6 h

and US\$70 using widely available commercial cloud nodes. This cost and turnaround time is much more amenable to rapid prototyping and parameter exploration than even the fastest competing method (Wtdbg2), which was on average 7.5 times slower and 3.7 times more expensive.

The combination of the Shasta assembler and nanopore long-read sequences produced using the PromethION sequencer realizes substantial improvements in throughput; we completed all 2.3 Tb of nanopore data collection in 9 d, running up to 15 flow cells simultaneously.

In terms of assembly, we obtained an average NG50 of 18.5 Mb for the 11 genomes, roughly three times higher than for the first nanopore-sequenced human genome, and comparable with the best achieved by alternative technologies^{13,49}. We found the addition of HiC sequencing for scaffolding necessary to achieve chromosome scale assemblies. However, our results are consistent with previous modeling based on the size and distribution of large repeats in the human genome, which predicts that an assembly based on 30 times coverage of such reads of >100 kb would approach the continuity of complete human chromosomes^{18,32}.

Relative to alternate long-read and linked-read sequencing, the read identity of nanopore reads is lower, however, improving over time^{12,18}. We observe modal read identity of 92.5%, resulting in better than QV = 30 base quality for haploid polished assembly from nanopore reads alone. The accurate resolution of highly repetitive and recently duplicated sequence will depend on long-read polishing, because short reads are generally not uniquely mappable. Further polishing using complementary data types, including PacBio HiFi reads⁴⁹ and 10X Chromium⁵⁰, will likely prove useful in achieving QV 40+ assemblies.

Shasta produces a notably more conservative assembly than competing tools, trading greater correctness for contiguity and total produced sequence. For example, the ratio of total length to aligned length is relatively constant for all other assemblers, where approximately 1.6% of sequence produced does not align across the three evaluated samples. In contrast, on average just 0.38% of Shasta's sequence does not align to GRCh38, representing a more than four times reduction in unaligned sequence. Additionally, we note substantially lower disagreement counts, resulting in much smaller differences between the raw NGx and corrected NGAx values. Shasta also produces substantially more base-level accurate assemblies than the other competing tools. MarginPolish and HELEN provide a consistent improvement of base quality over

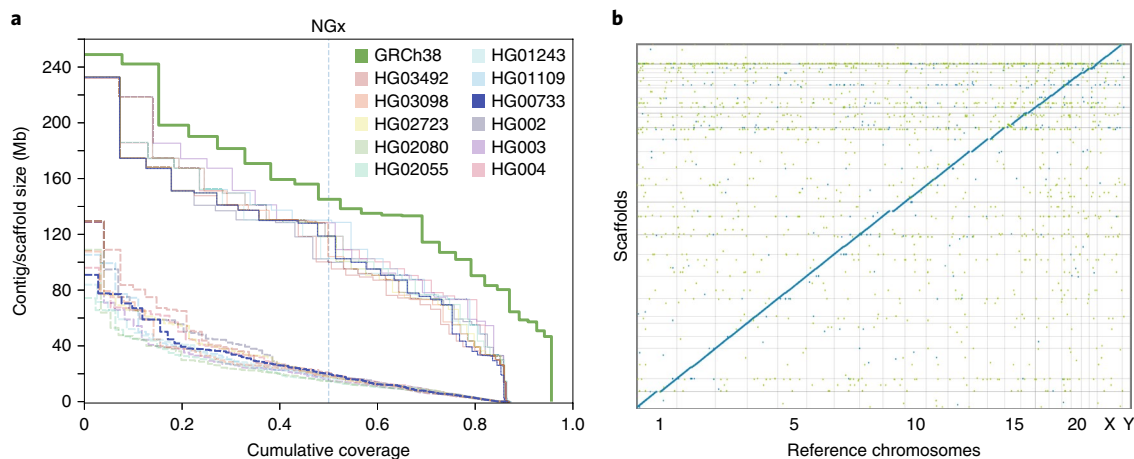


Fig. 5 | HiRise scaffolding for 11 genomes. a, NGx plots for each of the 11 genomes, before (dashed) and after (solid) scaffolding with HiC sequencing reads, GRCh38 minus alternate sequences is shown for comparison. **b,** Dot plot showing alignments between the scaffolded HG00733 Shasta assembly and GRCh38 chromosome scaffolds. Blue indicates forward aligning segments, green indicates reverse, with both indicating unique alignments.

all tested assemblers, with more accurate results than the current state-of-the-art long-read polishing workflow.

We assembled and compared haploid, trio-binned and diploid samples. Trio-binned samples show great promise for haplotype assembly, for example contiguously assembling an MHC haplogroup, but the halving of effective coverage resulted in ultimately less contiguous human assemblies with higher base-error rates than the related, chimeric diploid assembly. This can potentially be rectified by merging the haplotype assemblies to produce a pseudo-haplotype or increasing sequencing coverage. Indeed, the improvements in contiguity and base accuracy in CHM13 over the diploid samples illustrate what can be achieved with higher coverage of a haploid sample. We believe that one of the most promising directions for the assembly of diploid samples is the integration of phasing into the assembly algorithm itself, as pioneered by others^{17,51,52}. We anticipate that the new tools we have described here are suited for this next step: the Shasta framework is well placed for producing phased assemblies over structural variants, MarginPolish is built off of infrastructure designed to phase long reads² and the HELEN model could be improved to include haplotagged features for the identification of heterozygous sites.

Connected together, the tools we present enabled a polished assembly to be produced in around 24h and for roughly US\$180, against the fastest comparable combination of Wtdbg2, Racon and Medaka that costs 5.3 times more and is 4.3 times slower while producing measurably worse results in terms of disagreements, contiguity and base-level accuracy. Substantial further parallelism of polishing, the main time drain in our current pipeline, is easily possible.

We are working toward the goal of having a half-day turnaround of our complete computational pipeline. With real-time basecalling, a DNA-to-de novo assembly could conceivably be achieved in less than 96h. Such speed would enable screening of human genomes for abnormalities in difficult-to-sequence regions.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41587-020-0503-6>.

Received: 25 January 2020; Accepted: 26 March 2020;

References

- McKenna, A. et al. The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- Ebler, J., Haukness, M., Pesout, T., Marschall, T. & Paten, B. Haplotype-aware diplotyping from noisy long reads. *Genome Biol.* **20**, e116 (2019).
- Zook, J. M. et al. An open resource for accurately benchmarking small variant and reference calls. *Nat. Biotechnol.* **37**, 561–566 (2019).
- Poplin, R. et al. A universal snp and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
- Bradnam, K. R. et al. Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *Gigascience* **2**, 10 (2013).
- Alkan, C., Coe, B. P. & Eichler, E. E. Genome structural variation discovery and genotyping. *Nat. Rev. Genet.* **12**, 363–376 (2011).
- Kosugi, S. et al. Comprehensive evaluation of structural variation detection algorithms for whole genome sequencing. *Genome Biol.* **20**, 1–18 (2019).
- Chaisson, M. J. P. et al. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* **10**, 1–16 (2019).
- Belton, J. M. et al. Hi-C: a comprehensive technique to capture the conformation of genomes. *Methods* **58**, 268–276 (2012).
- Falconer, E. & Lansdorp, P. M. Strand-seq: a unifying tool for studies of chromosome segregation. *Semin. Cell Dev. Biol.* **24**, 643–652 (2013).
- Weisenfeld, N. I., Kumar, V., Shah, P., Church, D. M. & Jaffe, D. B. Direct determination of diploid genome sequences. *Genome Res.* **27**, 757–767 (2017).
- Jain, M. et al. Improved data analysis for the MinION nanopore sequencer. *Nat. Methods* **12**, 351–356 (2015).
- Eid, J. et al. Real-time DNA sequencing from single polymerase molecules. *Science* **323**, 133–138 (2009).
- Huddleston, J. et al. Discovery and genotyping of structural variation from long-read haploid genome sequence data. *Genome Res.* **27**, 677–685 (2017).
- Sedlazeck, F. J. et al. Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).
- Patterson, M. D. et al. WhatsHap: weighted haplotype assembly for future-generation sequencing reads. *J. Comput. Biol.* **22**, 498–509 (2015).
- Chin, C. S. et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* **13**, 1050–1054 (2016).
- Jain, M. et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* **36**, 338–345 (2018).
- Eichler, E. E., Clark, R. A. & She, X. An assessment of the sequence gaps: unfinished business in a finished human genome. *Nat. Rev. Genet.* **5**, 345–354 (2004).
- Fiddes, I. T. et al. Human-specific NOTCH2NL genes affect notch signaling and cortical neurogenesis. *Cell* **173**, 1356–1369.e22 (2018).
- Jain, M. et al. Linear assembly of a human centromere on the Y chromosome. *Nat. Biotechnol.* **36**, 321–323 (2018).
- Altshuler, D. M. et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
- Zook, J. M. et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
- Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
- ruanjue/smartdenovo: ultra-fast de novo assembler using long noisy reads (GitHub, 2020); <https://github.com/ruanjue/smartdenovo>
- Miller, J. R. et al. Aggressive assembly of pyrosequencing reads with mates. *Bioinformatics* **24**, 2818–2824 (2008).
- Broder, A. Z. On the resemblance and containment of documents. In *Proc. International Conference on Compression and Complexity of Sequences* 21–29 (IEEE, 1997); <https://doi.org/10.1109/sequen.1997.666900>
- Berlin, K. et al. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* **33**, 623–630 (2015).
- Ruan, J. & Li, H. Fast and accurate long-read assembly with wtdbg2. *Nat. Methods* **17**, 155–158 (2020).
- Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
- Koren, S. et al. Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
- Miga, K. H. et al. Telomere-to-telomere assembly of a complete human X chromosome. Preprint at *bioRxiv* <https://doi.org/10.1101/735928> (2019).
- Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile genome assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142–i150 (2018).
- Audano, P. A. et al. Characterizing the major structural variant alleles of the human genome. *Cell* **176**, 663–675.e19 (2019).
- Sudmant, P. H. et al. Global diversity, population stratification, and selection of human copy-number variation. *Science* (80-.). **349**, aab3761 (2015).
- Zook, J. M. et al. A robust benchmark for germline structural variant detection. Preprint at *bioRxiv* <https://doi.org/10.1101/664623> (2019).
- Brandt, D. Y. C. et al. Mapping bias overestimates reference allele frequencies at the HLA genes in the 1000 genomes project phase I data. *G3. Genes, Genomes, Genet.* **5**, 931–941 (2015).
- Turner, T. R. et al. Single molecule real-time DNA sequencing of HLA genes at ultra-high resolution from 126 International HLA and Immunogenetics Workshop cell lines. *HLA* **91**, 88–101 (2018).
- Koren, S. et al. De novo assembly of haplotype-resolved genomes with trio binning. *Nat. Biotechnol.* **36**, 1174–1182 (2018).
- Lee, C., Grasso, C. & Sharlow, M. F. Multiple sequence alignment using partial order graphs. *Bioinformatics* **18**, 452–464 (2002).
- Medsker, L. & Jain, D. L. (eds) *Recurrent Neural Networks: Design and Applications* (CRC Press, 2001).
- Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
- Walker, B. J. et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
- Fiddes, I. T. et al. Comparative Annotation Toolkit (CAT)-simultaneous clade and personal genome annotation. *Genome Res.* **28**, 1029–1038 (2018).
- Harrow, J. et al. GENCODE: the reference human genome annotation for the ENCODE project. *Genome Res.* **22**, 1760–1774 (2012).
- Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).

47. Vollger, M. R. et al. Improved assembly and variant detection of a haploid human genome using single-molecule, high-fidelity long reads. *Ann. Hum. Genet.* **84**, 125–140 (2020).
48. Putnam, N. H. et al. Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res.* **26**, 342–350 (2016).
49. Wenger, A. M. et al. Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
50. Ma, Z. S., Li, L., Ye, C., Peng, M. & Zhang, Y. P. Hybrid assembly of ultra-long Nanopore reads augmented with 10x-Genomics contigs: demonstrated with a human genome. *Genomics* **111**, 1896–1901 (2019).
51. Garg, S. et al. A graph-based approach to diploid genome assembly. *Bioinformatics* **34**, i105–i114 (2018).
52. Levy, S. et al. The diploid genome sequence of an individual human. *PLoS Biol.* **5**, e254 (2007).

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



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 license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license 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 license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020

Methods

Sample selection. The goal of sample selection was to select a set of individuals that collectively captured the maximum amount of weighted allelic diversity⁵³. To do this, we created a list of all low-passage lymphoblastoid cell lines that are part of a trio available from the 1,000 Genomes Project collection⁵⁴ (we selected trios to allow future addition of pedigree information, and low-passage line to minimize acquired variation). In some cases, we considered the union of parental alleles in the trios due to not having genotypes for the offspring. Let a weighted allele be a variant allele and its frequency in the 1,000 Genomes Project Phase 3 confidence variant set (VCF). We selected the first sample from our list that contained the largest sum of frequencies of weighted alleles, reasoning that this sample should have the largest expected fraction of variant alleles in common with any other randomly chosen sample. We then removed the variant alleles from this first sample from the set of variant alleles in consideration and repeated the process to pick the second sample, repeating the process recursively until we had selected seven samples. This set greedily, heuristically optimizes the maximum sum of weighted allele frequencies in our chosen sample subset. We also added the three Ashkenazim Trio samples and the Puerto Rican individual (HG00733). These four samples were added for the purposes of comparison with other studies that are using them²³.

Cell culture. Lymphoblastoid cultures for each individual were obtained from the Coriell Institute Cell Repository (coriell.org) and were cultured in RPMI 1640 supplemented with 15% fetal bovine serum (Life Technologies). The cells underwent a total of six passages (p3 + 3). After expansion, cells were collected by pelleting at 300g for 5 min. Cells were resuspended in 10 ml of PBS and a cell count was taken using a BiRad TC20 cell counter. Cells were aliquoted into 50 ml of conical tubes containing 50 million cells, pelleted as above and washed with 10 ml of PBS before a final pelleting after which the PBS was removed and the samples were flash frozen on dry ice and stored at -80°C until ready for further processing.

DNA extraction and size selection. We extracted high-molecular weight DNA using the Qiagen Puregene kit. We followed the standard protocol with some modifications. Briefly, we lysed the cells by adding 3 ml of Cell Lysis Solution per 10 million cells, followed by incubation at 37°C for up to 1 h. We performed mild shaking intermediately by hand and avoided vortexing. Once clear, we split the lysate into 3-ml aliquots and added 1 ml of protein precipitation solution to each of the tubes. This was followed by pulse vortexing three times for 5 s each time. We next spun this at 2,000g for 10 min. We added the supernatant from each tube to a new tube containing 3 ml of isopropanol, followed by 50 \times inversion. The high-molecular weight DNA precipitated and formed a dense thread-like jelly. We used a disposable inoculation loop to extract the DNA precipitate. We then dipped the DNA precipitate, while it was on the loop, into ice-cold 70% ethanol. After this, the DNA precipitate was added to a new tube containing 50–250 μl 1 \times TE buffer. The tubes were heated at 50°C for 2 h and then left at room temperature overnight to allow resuspension of the DNA. The DNA was then quantified using Qubit and NanoDrop.

We used the Circulomics Short-Read Eliminator kit to deplete short fragments from the DNA preparation. We size selected 10 μg of DNA using the Circulomics recommended protocol for each round of size selection.

Nanopore sequencing. We used the SQK-LSK109 kit and its recommended protocol for making sequencing libraries. We used 1 μg of input DNA per library. We prepared libraries at a 3 \times scale since we performed a nuclease flush on every flow cell, followed by the addition of a fresh library.

We used the standard PromethION scripts for sequencing. At around 24 h, we performed a nuclease flush using the ONT recommended protocol. We then reprimed the flow cell, and added a fresh library corresponding to the same sample. After the first nuclease flush, we restarted the run setting the voltage to -190 mV . We repeated the nuclease flush after another around 24 h (that is, around 48 h into sequencing), reprimed the flow cell, added a fresh library and restarted the run setting the run voltage to -200 mV .

We performed basecalling using Guppy v.2.3.5 on the PromethION tower using the graphics processing units (GPUs). We used the MinION DNA flipflop model (dna_r9.4.1_450bps_flipflop.cfg), as recommended by ONT.

Chromatin crosslinking and extraction from human cell lines. We thawed the frozen cell pellets and washed them twice with cold PBS before resuspension in the same buffer. We transferred aliquots containing five million cells by volume from these suspensions to separate microcentrifuge tubes before chromatin crosslinking by addition of paraformaldehyde (EMS catalog no. 15714) to a final concentration of 1%. We briefly vortexed the samples and allowed them to incubate at room temperature for 15 min. We pelleted the crosslinked cells and washed them twice with cold PBS before thoroughly resuspending in lysis buffer (50 mM Tris-HCl, 50 mM NaCl, 1 mM EDTA, 1% SDS) to extract crosslinked chromatin.

The HiC method. We bound the crosslinked chromatin samples to SPRI beads, washed three times with SPRI wash buffer (10 mM Tris-HCl, 50 mM NaCl, 0.05% Tween-20) and digested by DpnII (20 U, NEB catalog no. R0543S) for 1 h at 37°C

in an agitating thermal mixer. We washed the bead-bound samples again before incorporation of Biotin-11-dCTP (ChemCyte catalog no. CC-6002-1) by DNA Polymerase I, Klenow Fragment (10 U, NEB catalog no. M0210L) for 30 min at 25°C with shaking. Following another wash, we carried out blunt-end ligation by T4 DNA Ligase (4,000 U, NEB Catalog No. M0202T) with shaking overnight at 16°C . We reversed the chromatin crosslinks, digested the proteins, eluted the samples by incubation in crosslink reversal buffer (5 mM CaCl₂, 50 mM Tris-HCl, 8% SDS) with Proteinase K (30 μg , Qiagen catalog no. 19133) for 15 min at 55°C followed by 45 min at 68°C .

Sonication and Illumina library generation with biotin enrichment. After SPRI bead purification of the crosslink-reversed samples, we transferred DNA from each to Covaris microTUBE AFA Fiber Snap-Cap tubes (Covaris catalog no. 520045) and sonicated to an average length of 400 ± 85 base pairs using a Covaris ME220 Focused-Ultrasonicator. Temperature was held stably at 6°C and treatment lasted 65 s per sample with a peak power of 50 W, 10% duty factor and 200 cycles per burst. The average fragment length and distribution of sheared DNA was determined by capillary electrophoresis using an Agilent FragmentAnalyzer 5200 and HS NGS Fragment Kit (Agilent catalog no. DNF-474-0500). We ran sheared DNA samples twice through the NEBNext Ultra II DNA Library Prep Kit for Illumina (catalog no. E7645S) End Preparation and Adapter Ligation steps with custom Y adapters to produce library preparation replicates. We purified ligation products via SPRI beads before Biotin enrichment using Dynabeads MyOne Streptavidin C1 beads (ThermoFisher catalog no. 65002).

We performed indexing PCR on streptavidin beads using KAPA HiFi HotStart ReadyMix (catalog no. KK2602) and PCR products were isolated by SPRI bead purification. We quantified the libraries by Qubit 4 fluorometer and FragmentAnalyzer 5200 HS NGS Fragment Kit (Agilent catalog no. DNF-474-0500) before pooling for sequencing on an Illumina HiSeq X at Fulgent Genetics.

Analysis methods. Read alignment identities. To generate the identity violin plots (Fig. 1c,e) we aligned all the reads for each sample and flow cell to GRCh38 using minimap2 (ref. ²⁴) with the map-ont preset. Using a custom script `get_summary_stats.py` in the repository https://github.com/rlorigro/nanopore_assembly_and_polishing_assessment, we parsed the alignment for each read and enumerated the number of matched (N_m), mismatched (N_x), inserted (N_i) and deleted (N_d) bases. From this, we calculated alignment identity as $N_m/(N_m + N_x + N_i + N_d)$. These identities were aggregated over samples and plotted using the seaborn library with the script `plot_summary_stats.py` in the same repository. This method was used to generate Fig. 1c,e. For Fig. 1e, we selected reads from HG00733 flowcell1 aligned to GRCh38 chr1. The 'Standard' identities are used from the original reads/alignments. To generate identity data for the 'RLE' portion, we extracted the reads above, run-length encoded the reads and chr1 reference, and followed the alignment and identity calculation process described before. Sequences were run-length encoded using a simple script — https://github.com/rlorigro/runlength_analysis/blob/master/runlength_encode_fasta.py — and aligned with minimap2 using the map-ont preset and `-k 19`.

Base-level error-rate analysis with Pomoxis. We analyzed the base-level error rates of the assemblies using the `assess_assembly` tool of the Pomoxis toolkit (<https://github.com/nanoporetech/pomoxis>). The `assess_assembly` tool is tailored to compute the error rates in a given assembly compared to a truth assembly. It reports an identity error rate, insertion error rate, deletion error rate and an overall error rate. The identity error rate indicates the number of erroneous substitutions, the insertion error rate is the number of incorrect insertions and the deletion error rate is the number of deleted bases averaged over the total aligned length of the assembly to the truth. The overall error rate is the sum of the identity, insertion and deletion error rates. For the purpose of simplification, we used the indel error rate, which is the sum of insertion and deletion error rates.

The `assess_assembly` script takes an input assembly and a reference assembly to compare against. The assessment tool chunks the reference assembly to 1-kb regions and aligns it back to the input assembly to get a trimmed reference. Next, the input is aligned to the trimmed reference sequence with the same alignment parameters to get an input assembly to the reference assembly alignment. The total aligned length is the sum of the lengths of the trimmed reference segments where the input assembly has an alignment. The total aligned length is used as the denominator while averaging each of the error categories to limit the assessment in only correctly assembled regions. Then the tool uses `stats_from_bam`, which counts the number of mismatch bases, insert bases and delete bases at each of the aligned segments, and reports the error rate by averaging them over the total aligned length.

Truth assemblies for base-level error-rate analysis. We used HG002, HG00733 and CHM13 for base-level error-rate assessment of the assembler and the polisher. These three assemblies have high-quality assemblies publicly available, which are used as the ground truth for comparison. Two of the samples, HG002 and HG00733, are diploid samples; hence, we picked one of the two possible haplotypes as the truth. The reported error rate of HG002 and HG00733 include some errors arising due to the zygosity of the samples. The complete hydatidiform mole sample

CHM13 is a haploid human genome that is used to assess the applicability of the tools on haploid samples. We have gathered and uploaded all the files we used for assessment in one place at https://console.cloud.google.com/storage/browser/kishwar-helen/truth_assemblies/.

To generate the HG002 truth assembly, we gathered the publicly available GIAB high-VCF against GRCh38 reference sequence. Then we used bedtools to create an assembly (FASTA) file from the GRCh38 reference and the high-confidence variant set. We got two files using this process for each of the haplotypes, and we picked one randomly as the truth. All the diploid HG002 assembly is compared against this one chosen assembly. GIAB also provides a bed file annotating high-confidence region where the called variants are highly precise and sensitive. We used this bed file with `assess_assembly` to ensure that we compare the assemblies only in the high-confidence regions.

The HG00733 truth is from the publicly available phased PacBio high-quality assembly of this sample⁵⁵. We picked phase0 as the truth assembly and acquired it from the National Center for Biotechnology Information under accession [GCA_003634895.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_003634895.1). We note that the assembly is phased but not haplotyped, such that portions of phase0 will include sequences from both parental haplotypes and is not suitable for trio-binned analyses. Furthermore, not all regions were fully phased; regions with variants that are represented as some combination of both haplotypes will result in lower QV and a less accurate truth.

For CHM13, we used the v.0.6 release of CHM13 assembly by the T2T consortium³². The reported quality of this truth assembly in Phred quality score (QV) value is 39. One of the attributes of this assembly is chromosome X. As reported by the T2T assembly authors, chromosome X of CHM13 is the most complete (end-to-end) and high-quality assembly of any human chromosome. We obtained the chromosome X assembly, which is the highest-quality truth assembly (QV \geq 40) we have.

QUAST/BUSCO. To quantify contiguity, we primarily depended on the tool QUAST³³. QUAST identifies misassemblies as main rearrangement events in the assembly relative to the reference. We use the phrase ‘disagreement’ in our analysis, as we find ‘misassembly’ inappropriate considering potentially true structural variation. For our assemblies, we quantified all contiguity stats against GRCh38, using autosomes plus chromosomes X and Y only. We report the total disagreements given that their relevant ‘size’ descriptor was greater than 1 kb, as is the default behavior in QUAST. QUAST provides other contiguity statistics in addition to disagreement count, notably total length and total aligned length as reported in Fig. 2d. To determine total aligned length (and unaligned length), QUAST performs collinear chaining on each assembled contig to find the best set of nonoverlapping alignments spanning the contig. This process contributes to QUAST’s disagreement determination. We consider an unaligned sequence to be the portions of the assembled contigs that are not part of this best set of nonoverlapping alignments. All statistics are recorded in Supplementary Table 5. For all QUAST analyses, we used the flags `min-identity 80` and `fragmented`.

QUAST also produces an NGAx plot (similar to an NGx plot) that shows the aligned segment size distribution of the assembly after accounting for disagreements and unalignable regions. The intermediate segment lengths that would allow NGAx plots to be reproduced across multiple samples on the same axis (as is shown in Fig. 2b) are not stored, so we created a GitHub fork of QUAST to store this data during execution at <https://github.com/rlogigro/quast>. Finally, the assemblies and the output of QUAST were parsed to generate figures with an NGx visualization script, `ngx_plot.py`, found at http://github.com/rlogigro/nanopore_assembly_and_polishing_assessment/.

For NGx and NGAx plots, a total genome size of 3.23 Gb was used to calculate cumulative coverages.

BUSCO⁴⁶ is a tool that quantifies the number of benchmarking universal single-copy orthologs present in an assembly. We ran BUSCO via the option within QUAST, comparing against the eukaryota set of orthologs from OrthoDB v.9.

Disagreement assessments. To analyze the QUAST-reported disagreements for different regions of the genome, we gathered the known segmental duplication regions⁸, centromeric regions for GRCh38 and known regions in GRCh38 with structural variation for HG002 from GIAB³⁶. We used a Python script `quast_sv_extractor.py` that compares each reported disagreement from QUAST to the segmental duplication, SV and centromeric regions and discounts any disagreement that overlaps with these regions. The `quast_sv_extractor.py` script can be found at <https://github.com/kishwarshafin/helen/tree/master/helen/modules/python/helper>.

The segmental duplication regions of GRCh38 defined in the `ucsc.collapsed.sorted.segdup` file can be downloaded from <https://github.com/mvollger/segDupPlots/>.

The defined centromeric regions of GRCh38 for all chromosomes are used from the available summary at <https://www.ncbi.nlm.nih.gov/grc/human>.

For GIAB HG002, known SVs for GRCh38 are available in NIST SVs Integration v.0.6/ under `ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/`. We used the Tier1+2 bed file available at the GIAB ftp site.

We further exclude SV enriched regions like centromeres, secondary constriction regions, acrocentric arms, large tandem repeat arrays, segmental duplications and the Y chromosome plus 10-kb pairs on either side of them. The file is available at https://github.com/kishwarshafin/helen/blob/master/masked_regions/GRCh38_masked_regions.bed.

To analyze disagreements within the intersection of the assembled sequences we performed the following analysis. For each assembly we used `minimap2` and `samtools` to create regions of unique alignment to GRCh38. For `minimap2` we used the options `--secondary=no -a --eqx -Y -x asm20 -m 10000 -z 10000,50 -r 50000 --end-bonus=100 -O 5,6 -E 4,1 -B 5`. We fed these alignments into `samtools view` with options `-F 260 -u` and then `samtools sort` with option `-m`. We then scanned 100 basepair windows of GRCh38 to find windows where all assemblies for the given sample were aligned with a one-to-one mapping to GRCh38. We then report the sum of disagreements across these windows. The script for this analysis can be found at https://github.com/mvollger/consensus_regions.

Trio-binning. We performed trio-binning on two samples, HG002 and HG00733 (ref. ³⁹). For HG00733, we obtained the parental read sample accessions (HG00731, HG00732) from the 1,000 Genome Database. Then we counted *k*-mers with `meryl` to create maternal and paternal *k*-mer sets. Based on manual examination of the *k*-mer count histograms to determine an appropriate threshold, we excluded *k*-mers occurring fewer than six times for maternal set and five times for paternal set. We subtracted the paternal set from the maternal set to get *k*-mers unique to the maternal sample and similarly derived unique paternal *k*-mer set. Then for each read, we counted the number of occurrences of unique maternal and paternal *k*-mers and classified the read based on the highest occurrence count. During classification, we avoided normalization by *k*-mer set size. This resulted in 35.2 \times maternal, 37.3 \times paternal and 5.6 \times unclassified for HG00733. For HG002, we used the Illumina data for the parental samples (HG003, HG004) from GIAB project³³. We counted *k*-mers using `meryl` and derived maternal paternal sets using the same protocol. We filtered *k*-mers that occur fewer than 25 times in both maternal and paternal sets. The classification resulted in 24 \times maternal, 23 \times paternal and 3.5 \times unknown. The commands and data source are detailed in the Supplementary Information.

Transcript analysis with comparative annotation toolkit (CAT). We ran the CAT⁴⁴ to annotate the polished assemblies to analyze how well Shasta assembles transcripts and genes. Each assembly was individually aligned to the GRCh38 reference assembly using `Cactus`⁵⁶ to create the input alignment to CAT. The GENCODE⁴⁵ V30 annotation was used as the input gene set. CAT was run in the `transMap` mode only, without Augustus refinement, since the goal was only to evaluate the quality of the projected transcripts. All transcripts on chromosome Y were excluded from the analysis since some samples lacked a Y chromosome.

Run-length confusion matrix. To generate run-length confusion matrices from reads and assemblies, we run-length encoded the assembly/read sequences and reference sequences using a purpose-built Python script, `measure_runlength_distribution_from_fasta.py`. The script requires a reference and sequence file and can be found in the GitHub repository https://github.com/rlogigro/runlength_analysis/. The run-length encoded nucleotides were aligned to the run-length encoded reference nucleotides with `minimap2`. As run-length encoded sequences cannot have identical adjacent nucleotides, the number of unique *k*-mers is diminished with respect to standard sequences. As `minimap2` uses empirically determined sizes for seed *k*-mers, we used a *k*-mer size of 19 to approximately match the frequency of the default size (15) used by the presets for standard sequences. For alignment of reads and assemblies we used the `map-ont` and `asm20` presets, respectively.

By iterating through the alignments, each match position in the cigar string (mismatched nucleotides are discarded) was used to find a pair of lengths (*x*, *y*) such that *x* is a predicted length and *y* is the true (reference) length. For each pair, we updated a matrix that contains the frequency of every possible pairing of prediction versus truth, from length 1 to 50 bp. Finally, this matrix is normalized by dividing each element by the sum of the observations for its true run length and plotted as a heatmap. Each value represents the probability of predicting a length for a given true length.

Runtime and cost analysis. Our runtime analysis was generated with multiple methods detailing the amount of time the processes took to complete. These methods include the Unix command time and a home-grown resource tracking script, which can be found in the <https://github.com/rlogigro/TaskManager> repository. We note that the assembly and polishing methods have different resource requirements, and do not all fully use available CPUs, GPUs and memory over the program’s execution. As such, we report runtimes using wall-clock time and the number of CPUs the application was configured to use, but do not convert to CPU hours. Costs reported in the figures are the product of the runtime and AWS instance price. Because portions of some applications do not fully use CPUs, cost could potentially be reduced by running on a smaller instance that would be fully used, and runtime could be reduced by running on a larger instance that can be fully used for some portion of execution. We particularly note the long runtime

of Medaka and found that for most of the total runtime, only a single CPU was used. Last, we note that data transfer times are not reported in runtimes. Some of the data required or generated exceeds hundreds of gigabytes, which could be notable in relation to the runtime of the process. Notably, the images generated by MarginPolish and consumed by HELEN were often greater than 500 GB in total.

All recorded runtimes are reported in the supplement. For Shasta, times were recorded to the tenth of the hour. All other runtimes were recorded to the minute. All runtimes reported in figures were run on the AWS cloud platform.

Shasta runtime reported in Fig. 2f was determined by averaging across all 12 samples. Wtdbg2 runtime was determined by summing runtimes for wtdbg2 and wtpoa-cns and averaging across the HG00733, HG002 and CHM13 runs. Flye runtime was determined by averaging across the HG00733, HG002 and CHM13 runs, which were performed on multiple instance types (x1.16xlarge and x1.32xlarge). We calculated the total cost and runtime for each run and averaged these amounts; no attempt to convert these to a single instance type was performed. Precise Canu runtimes are not reported, as they were run on the National Institutes of Health (NIH) Biowulf cluster. Each run was restricted to nodes with 28 cores (56 hyperthreads) (2×2680v4 or 2×2695v3 Intel CPUs) and 248 GB of RAM or 16 cores (32 hyperthreads) (2×2650v2 Intel CPUs) and 121 GB of RAM. Full details of the cluster are available at <https://hpc.nih.gov>. The runs took between 219,000 and 223,000 CPU hours (4–5 wall-clock days). No single job used more than 80 GB of RAM/12 CPUs. We find the r5.xlarge (US\$1,008 per hour) to be the cheapest AWS instance type possible considering this resource usage, which puts estimated cost between US\$18,000 and US\$19,000 per genome.

For MarginPolish, we recorded all runtimes, but used various thread counts that did not always fully use the instance's CPUs. The runtime reported in the figure was generated by averaging across eight of the 12 samples, selecting runs that used 70 CPUs (of the 72 available on the instance). The samples this was true for were GM24385, HG03492, HG01109, HG02055, HG02080, HG01243, HG03098 and CHM13. Runtimes for read alignments used by MarginPolish were not recorded. Because MarginPolish requires an aligned BAM, we found it unfair to not report this time in the figure as it is a required step in the workflows for MarginPolish, Racon and Medaka. As a proxy for the unrecorded read alignment time used to generate BAMs for MarginPolish, we added the average alignment time recorded while aligning reads in preparation for Medaka runs. We note that the alignment for MarginPolish was done by piping output from minimap2 directly into samtools sort, and piping this into samtools view to filter for primary and supplementary reads. Alignment for Medaka was done using mini_align, which is a wrapper for minimap2 bundled in Medaka that simultaneously sorts output.

Reported HELEN runs were performed on GCP except for HG03098, but on instances that match the AWS instance type p2.8xlarge in both CPU count and GPU (NVIDIA Tesla P100). As such, the differences in runtime between the platforms should be negligible, and we have calculated cost based on the AWS instance price for consistency. The reported runtime is the sum of time taken by call_consensus.py and stitch.py. Unannotated runs were performed on UCSC hardware. Racon runtimes reflect the sum of four series of read alignment and polishing. The time reported in the figure is the average of the runtime of this process run on the Shasta assembly for HG00733, HG002 and CHM13.

Medaka runtime was determined by averaging across the HG00733, HG002 and CHM13 runs after running Racon four times on the Shasta assembly. We again note that this application in particular did not fully use the CPUs for most of the execution, and in the case of HG00733 appeared to hang and was restarted. The plot includes the average runtime from read alignment using minialign; this is separated in the tables in the Supplementary Information. We ran Medaka on an x1.16xlarge instance, which had more memory than was necessary. When determining cost, we chose to price the run based on the cheapest AWS instance type that we could have used accounting for configured CPU count and peak memory usage (c5n.18xlarge). This instance could have supported eight more concurrent threads, but as the application did not fully use the CPUs, we find this to be a fair representation.

Assembly of MHC. Each of the eight GRCh38 MHC haplotypes were aligned using minimap2 (with preset asm20) to whole-genome assemblies to identify spanning contigs. These contigs were then extracted from the genomic assembly and used for alignment visualization. For dot plots, Nucmer 4.0⁵⁷ was used to align each assembler's spanning contigs to the standard chr6:28000000-34000000 MHC region, which includes 500-Mb flanks. Output from this alignment was parsed with Dot⁵⁸, which has a web-based graphical user interface for visualization. All defaults were used in both generating the input files and drawing the figures. Coverage plots were generated from reads aligned to chr6, using a script, find_coverage.py, located at http://github.com/rlorigro/nanopore_assembly_and_polishing_assessment/.

The best matching alt haplotype (to Shasta, Canu and Flye) was chosen as a reference haplotype for quantitative analysis. Haplotypes with the fewest supplementary alignments across assemblers were top candidates for QUAST analysis. Candidates with comparable alignments were differentiated by identity. The highest contiguity/identity MHC haplotype was then analyzed with QUAST using -min-identity 80. For all MHC analyses regarding Flye, the unpolished output was used.

BAC analysis. At a high level, the BAC analysis was performed by aligning BACs to each assembly, quantifying their resolution and calculating identity statistics on those that were fully resolved.

We obtained 341 BACs for CHM13 (refs. ^{59,60}) and 179 for HG00733 (ref. ⁸) (complete BAC clones of VMRC62), which had been selected primarily by targeting complex or highly duplicated regions. We performed the following analysis on the full set of BACs (for CHM13 and HG00733), and a subset selected to fall within unique regions of the genome. To determine this subset, we selected all BACs that are greater than 10 kb away from any segmental duplication, resulting in 16 of HG00733 and 31 of CHM13. This subset represents simple regions of the genome that we would expect all assemblers to resolve.

For the analysis, BACs were aligned to each assembly with the command `minimap2 -secondary=no -t 16 -ax asm20 assembly.fasta bac.fasta>assembly.sam` and converted to a PAF-like format that describes aligned regions of the BACs and assemblies. Using this, we calculated two metrics describing how resolved each BAC was: closed is defined as having 99.5% of the BAC aligned to a single locus in the assembly; attempted is defined as having a set of alignments covering ≥95% of the BAC to a single assembly contig where all alignments are at least 1 kb away from the contig end. If such a set exists, it counts as attempted. We further calculated median and mean identities (using the alignment identity metric described above) of the closed BACs. These definitions were created such that a contig that is counted as attempted but not closed likely reflects a disagreement. The code for this can be found at <https://github.com/skoren/bacValidation>.

Short-read polishing. Chromosome X of the CHM13 assembly (assembled first with Shasta, then polished with MarginPolish and HELEN) was obtained by aligning the assembly to GRCh38 (using minimap2 with the -x asm20 flag). The 10X Chromium reads were downloaded from the Nanopore Whole Genome Sequencing Consortium (<https://github.com/nanopore-wgs-consortium/CHM13/>).

The 10X reads were from a NovaSeq eitinstrument at a coverage of approximately 50×. The reads corresponding to chromosome X were extracted by aligning the entire read set to the whole CHM13 assembly using the 10X Genomics Long Ranger Align pipeline (v.2.2), then extracting those corresponding to the corresponding chromosome X contigs with samtools. Pilon⁶³ was run iteratively for a total of three rounds, in each round aligning the reads to the current assembly with Long Ranger and then running Pilon with default parameters.

Structural variant assessment. To create an assembly graph in GFA format, Shasta v.0.1.0 was run using the HG002 sequence data with -MarkerGraph, simplifyMaxLength 10 to reduce bubble removal and -MarkerGraph, highCoverageThreshold 10 to reduce the removal of edges normally removed by the transitive reduction step.

To detect structural variation inside the assembly graphs produced by Shasta, we extracted unitigs from the graph and aligned them back to the linear reference. Unitigs are walks through the assembly graph that do not traverse any node end that includes a bifurcation. We first processed the Shasta assembly graphs with gimbricate (<https://github.com/ekg/gimbricate>). We used gimbricate to recompute overlaps in nonRLE space and to remove nodes in the graph only supported by a single sequencing read.

To remove overlaps from the graph edges, we then 'bluntified' resulting GFAs with `vg find -F` (<https://github.com/vgteam/vg>). We then applied `odgi unitig` (<https://github.com/vgteam/odgi>) to extract unitigs from the graph, with the condition that the starting node in the unitig generation must be at least 100 bp long. To ensure that the unitigs could be mapped back to the linear reference, we appended a random walk of 25 kb after the natural end of each unitig, with the expectation that even should unitigs would yield around 50 kb of mappable sequence. Finally, we mapped the unitigs to GRCh38 with minimap2 with a bandwidth of 25 kb (-r25000) and called variants in the alignments using `paftools.js` from the minimap2 distribution. We implemented the process in a single script that produces variant calls from the unitig set of a given graph: https://github.com/ekg/shastaGFA/blob/master/shastaGFAtoVCF_unitig_paftools.sh.

The extracted variants were compared to the structural variants from the GIAB benchmark in HG002 (v.0.6, ref. ³⁶). Precision, recall and F1 scores were computed on variants not overlapping simple repeats and within the benchmark's high-confidence regions. Deletions in the assembly and the GIAB benchmark were matched if they had at least 50% reciprocal overlap. Insertions were matched if located at less than 100 bp from each other and similar in size (50% reciprocal similarity).

Shasta. The following describes Shasta v.0.1.0 (<https://github.com/chanzuckerberg/shasta/releases/tag/0.1.0>), which was used throughout our analysis. All runs were done on an AWS x1.32xlarge instance (1,952 GB memory, 128 virtual processors). The runs used the Shasta recommended options for best performance (-memoryMode filesystem -memoryBacking 2M). Rather than using the distributed version of the release, the source code was rebuilt locally for best performance as recommended by Shasta documentation.

RLE of input reads. Shasta represents input reads using RLE. The sequence of each input read is represented as a sequence of bases, each with a repeat count that says

how many times each of the bases is repeated. Such a representation has previously been used in biological sequence analysis^{24–26}.

For example, the read
CGATTTAAGTTA
is represented as follows using RLE:
CGATAGTA
11132121

Using RLE makes the assembly process less sensitive to errors in the length of homopolymer runs, which are the most common type of errors in Oxford Nanopore reads. For example, consider these two reads:

CGATTTAAGTTA
CGATTAAGGGTTA
Using their raw representation, these reads can be aligned like this:
CGATTTAAG--TTA
CGATT--AAGGGTTA

Aligning the second read to the first required a deletion and two insertions. But in RLE, the two reads become:

CGATAGTA
11132121
CGATAGTA
11122321

The sequence portions are now identical and can be aligned trivially and exactly, without any insertions or deletions:

CGATAGTA
CGATAGTA

The differences between the two reads only appear in the repeat counts:
11132121
11122321

The Shasta assembler uses 1 byte to represent repeat counts, and as a result it only represents repeat counts between 1 and 255. If a read contains more than 255 consecutive bases, it is discarded on input. In the data we have analyzed so far, such reads are extremely rare.

Some properties of base sequences in RLE.

- In the sequence portion of the RLE, consecutive bases are always distinct. If they were not, the second one would be removed from the RLE sequence, while increasing the repeat count for the first one.
- With ordinary base sequences, the number of distinct k -mers of length k is 4^k . But with run-length base sequences, the number of distinct k -mers of length k is $4 \times 3^{k-1}$. This is a consequence of the previous bullet.
- The run-length sequence is generally shorter than the raw sequence and cannot be longer. For a long random sequence, the number of bases in the run-length representation is three-quarters of the number of bases in the raw representation.

Markers. Even with RLE, errors in input reads are still frequent. To further reduce sensitivity to errors, and also to speed up some of the computational steps in the assembly process, the Shasta assembler also uses a read representation based on markers. Markers are occurrences in reads of a predetermined subset of short k -mers. By default, Shasta uses for this purpose k -mers with $k = 10$ in RLE, corresponding to an average approximately 13 bases in raw read representation.

Just for the purposes of illustration, consider a description using markers of length 3 in RLE. There is a total $4 \times 3^2 = 36$ distinct such markers. We arbitrarily choose the following fixed subset of the 36, and we assign an identity to each of the k -mers in the subset as follows:

TGC 0
GCA 1
GAC 2
CGC 3

Consider now the following portion of a read in run-length representation (here, the repeat counts are irrelevant and so they are omitted):

CGACACGTATGCGCACGCTGCGCTCTGCAGC
GAC TGC CGC TGC
CGC TGC GCA
GCA CGC

Occurrences of the k -mers defined in the example above are shown and define the markers in this read. Note that markers can overlap. Using the marker identities defined in the list, we can summarize the sequence of this read portion as follows:

2 0 3 1 3 0 3 0 1

This is the marker representation of this read portion. It just includes the sequence of markers occurring in the read, not their positions. Note that the marker representation loses information, as it is not possible to reconstruct the complete initial sequence from the marker representation. This also means that the marker representation is insensitive to errors in the sequence portions that do not belong to any markers.

The Shasta assembler uses a random choice of the k -mers to be used as markers. The length of the markers k is controlled by assembly parameter `Kmers.k` with a default value of ten. Each k -mer is randomly chosen to be used as a marker

with probability determined by assembly parameter `--Kmers.probability` with a default value of 0.1. With these default values, the total number of distinct markers is approximately $0.1 \times 4 \times 3^9 \cong 7,900$.

The only constraint used in selecting k -mers to be used as markers is that if a k -mer is a marker, its reverse complement should also be a marker. This makes it easy to construct the marker representation of the reverse complement of a read from the marker representation of the original read. It also ensures strand symmetry in some of the computational steps.

It is possible that the random selection of markers is not optimal, and that it may be best to select the markers based on their frequency in the input reads or other criteria. These possibilities have not yet been investigated. Extended Data Fig. 1 shows the run-length representation of a portion of a read and its markers, as displayed by the Shasta http server.

Marker alignments. The marker representation of a read is a sequence in an alphabet consisting of the marker identities. This sequence is much shorter than the original sequence of the read but uses a much larger alphabet. For example, with default Shasta assembly parameters, the marker representation is ten times shorter than the run-length encoded read sequence, or about 13 times shorter than the raw read sequence. Its alphabet has around 8,000 symbols, many more than the four symbols that the original read sequence uses.

Because the marker representation of a read is a sequence, we can compute an alignment of two reads directly in marker representation. Computing an alignment in this way has two important advantages:

- The shorter sequences and larger alphabet make the alignment much faster to compute.
- The alignment is insensitive to read errors in the portions that are not covered by any marker.

For these reasons, the marker representation is orders of magnitude more efficient than the raw base representation when computing read alignments. Extended Data Fig. 2 shows an example alignment matrix.

Computing optimal alignments in marker representation. To compute the (likely) optimal alignment (example highlighted in green in Extended Data Fig. 2), the Shasta assembler uses a simple alignment algorithm on the marker representations of the two reads to be aligned. It effectively constructs an optimal path in the alignment matrix, but using some ‘banding’ heuristics to speed up the computation:

- The maximum number of markers that an alignment can skip on either read is limited to a maximum, under control of assembly parameter `Align.maxSkip` (default value 30 markers, corresponding to around 400 bases when all other Shasta parameters are at their default). This reflects the fact that Oxford Nanopore reads can often have long stretches in error. In the alignment matrix shown in Extended Data Fig. 2, there is a skip of about 20 markers (two light-gray squares) following the first ten aligned markers (green dots) on the top left.
- The maximum number of markers that an alignment can skip at the beginning or end of a read is limited to a maximum, under control of assembly parameter `Align.maxTrim` (default value 30 markers, corresponding to around 400 bases when all other Shasta parameters are at their default). This reflects the fact that Oxford Nanopore reads often have an initial or final portion that is not usable. These first two heuristics are equivalent to computing a reduced band of the alignment matrix.
- To avoid alignment artifacts, marker k -mers that are too frequent in either of the two reads being aligned are not used in the alignment computation. For this purpose, the Shasta assembler uses a criterion based on absolute number of occurrences of marker k -mers in the two reads, although a relative criterion (occurrences per kilobase) may be more appropriate. The current absolute frequency threshold is under control of assembly parameter `Align.maxMarkerFrequency` (default ten occurrences).

Using these techniques and with the default assembly parameters, the time to compute an optimal alignment is 10^{-3} – 10^{-2} s in the Shasta implementation as of release v.0.1.0 (April 2019). A typical human assembly needs to compute 10^8 read alignments that results in a total compute time $\sim 10^3$ – 10^6 s or $\sim 10^3$ – 10^4 s of elapsed time (1–3 h) on a machine with 128 virtual processors. This is one of the most computationally expensive portions of a Shasta assembly. Some additional optimizations are possible in the code that implement this computation and may be implemented in future releases.

Finding overlapping reads. Even though computing read alignments in marker representation is fast, it still is not feasible to compute alignments among all possible pairs of reads. For a human size genome with $\sim 10^6$ – 10^7 reads, the number of pairs to consider would be $\sim 10^{12}$ – 10^{14} , and even at 10^{-3} s per alignment the compute time would be $\sim 10^9$ – 10^{11} s or $\sim 10^7$ – 10^9 s elapsed time ($\sim 10^2$ – 10^4 d) when using 128 virtual processors.

Therefore, some means of narrowing down substantially the number of pairs to be considered is essential. The Shasta assembler uses for this purpose a slightly modified `MinHash`^{27,28} scheme based on the marker representation of reads.

In overview, the MinHash algorithm takes as input a set of items each characterized by a set of features. Its goal is to find pairs of the input items that have a high Jaccard similarity index: that is, pairs of items that have many features in common. The algorithm proceeds by iterations. At each iteration, a new hash table is created and a hash function that operates on the feature set is selected. For each item, the hash function of each of its features is evaluated, and the minimum hash function value found is used to select the hash table bucket that each item is stored in. It can be proved that the probability of two items ending up in the same bucket equals the Jaccard similarity index of the two items: that is, items in the same bucket are more likely to be highly similar than items in different buckets⁶¹. The algorithm then adds to the pairs of potentially similar items all pairs of items that are in the same bucket.

When all iterations are complete, the probability that a pair of items was found at least once is an increasing function of the Jaccard similarity of the two items. In other words, the pairs found are enriched for pairs that have high similarity. One can now consider all the pairs found (hopefully a much smaller set than all possible pairs) and compute the Jaccard similarity index for each, then keep only the pairs for which the index is sufficiently high. The algorithm does not guarantee that all pairs with high similarity will be found, only that the probability of finding all pairs is an increasing function of their similarity.

The algorithm is used by Shasta with items being oriented reads (a read in either original or reverse complemented orientation) and features being consecutive occurrences of m markers in the marker representation of the oriented read. For example, consider an oriented read with the following marker representation: 18,45,71,3,15,6,21

If m is selected as equal to four (the Shasta default, controlled by assembly parameter MinHash.m), the oriented read is assigned the following features:

(18,45,71,3)
(45,71,3,15)
(71,3,15,6)
(3,15,6,21)

From this picture of an alignment matrix in marker representation, we see that streaks of four or more common consecutive markers are relatively common. We have to keep in mind that, with Shasta default parameters, four consecutive markers span an average 40 bases in RLE or about 52 bases in the original raw base representation. At a typical error rate around 10%, such a portion of a read would contain on average five errors. Yet, the marker representation in run-length space is sufficiently robust that these common 'features' are relatively common despite the high error rate. This indicates that we can expect the MinHash algorithm to be effective in finding pairs of overlapping reads.

However, the MinHash algorithm has a feature that is undesirable for our purposes: namely, that the algorithm is good at finding read pairs with high Jaccard similarity index. For two sets X and Y , the Jaccard similarity index is defined as the ratio

$$J = \frac{X \cap Y}{X \cup Y}$$

Because the read length distribution of Oxford Nanopore reads is very wide, it is very common to have pairs of reads with very different lengths. Consider now two reads with lengths n_x and n_y , with $n_x < n_y$, that overlap exactly over the entire length n_x . The Jaccard similarity is in this case given by $n_x/n_y < 1$. This means that, if one of the reads in a pair is much shorter than the other one, their Jaccard similarity will be low even in the best case of exact overlap. As a result, the unmodified MinHash algorithm will not do a good job at finding overlapping pairs of reads with very different lengths.

For this reason, the Shasta assembler uses a small modification to the MinHash algorithm: instead of just using the minimum hash for each oriented read for each iteration, it keeps all hashes below a given threshold (this is not the same as keeping a fixed number of the lowest hashes for each read). Each oriented read can be stored in multiple buckets, one for each low hash encountered. The average number of low hashes on a read is proportional to its length, and, therefore, this change has the effect of eliminating the bias against pairs in which one read is much shorter than the other. The probability of finding a given pair is no longer driven by the Jaccard similarity. The modified algorithm is referred to as LowHash in the Shasta source code. Note that it is effectively equivalent to an indexing approach in which we index all features with low hash.

The LowHash algorithm is controlled by the following assembly parameters:

- MinHash.m (default 4): the number of consecutive markers that define a feature
- MinHash.hashFraction (default 0.01): the fraction of hash values that count as 'low'
- MinHash.minHashIterationCount (default 10): the number of iterations
- MinHash.maxBucketSize (default 10): the maximum number of items for a bucket to be considered. Buckets with more than this number of items are ignored. The goal of this parameter is to mitigate the effect of common repeats, which can result in buckets containing large numbers of unrelated oriented reads
- MinHash.minFrequency (default 2): the number of times a pair of oriented reads has to be found to be considered and stored as a possible pair of overlapping reads

Initial assembly steps. Initial steps of a Shasta assembly proceed as follows.

If the assembly is setup for best performance (--memoryMode filesystem --memoryBacking 2M if using the Shasta executable), all data structures are stored in memory and no disk activity takes place except for initial loading of the input reads, storing of assembly results and storing a small number of small files with useful summary information.

- (1) Input reads are read from FASTA files and converted to run-length representation.
- (2) k -mers to be used as markers are randomly selected.
- (3) Occurrences of those marker k -mers in all oriented reads are found.
- (4) The LowHash algorithm finds candidate pairs of overlapping oriented reads.
- (5) A marker alignment is computed for each candidate pair of oriented reads. If the marker alignment contains a minimum number of aligned markers, the pair is stored as an aligned pair. The minimum number of aligned markers is controlled by assembly parameter Align.minAlignedMarkerCount.

Read graph. Using the methods covered so far, an assembly has created a list of pairs of oriented reads, each pair having a plausible marker alignment. How to use this type of information for assembly is a classical problem with a standard solution⁶², the string graph.

It may be possible to adapt the prescriptions in the Myers paper to our situation in which a marker representation is used. However, we have not attempted this here, leaving it for future work.

Instead, the approach currently used in the Shasta assembler is very simple and can likely be improved. In the current simple approach, the Shasta assembler creates an undirected graph, the Read Graph, in which each vertex represents an oriented read (that is, a read in either original orientation or reverse complemented) and an undirected edge between two vertices is created if we have found an alignment between the corresponding oriented reads.

However, the read graph as constructed in this way suffers from high connectivity in repeat regions. Therefore, the Shasta assembler only keeps a k -nearest-neighbor subset of the edges. That is, for each vertex (oriented read) we only keep the k edges with the best alignments (greatest number of aligned markers). The number of edges kept for each vertex is controlled by assembly parameter ReadGraph.maxAlignmentCount, with a default value of six. Note that, despite the k -nearest-neighbor subset, it remains possible for a vertex to have degrees more than k .

Note that each read contributes two vertices to the read graph, one in its original orientation, and one in reverse complemented orientation. Therefore, the read graph contains two strands; each strand at full coverage. This makes it easy to investigate and potentially detect erroneous strand jumps that would be much less obvious if using approaches with one vertex per read.

An example of one strand is shown in Extended Data Fig. 3a. Even though the graph is undirected, edges that correspond to overlap alignments are drawn with an arrow that points from the prefix oriented read to the suffix one, to represent the direction of overlap. Edges that correspond to containment alignments (an alignment that covers one of the two reads entirely) are drawn in red and without an arrow. Vertices are drawn with area proportional to the length of the corresponding reads.

The linear structure of the read graph successfully reflects the linear arrangement of the input reads and their origin on the genome being assembled. However, deviations from the linear structure can occur in the presence of long repeats (Extended Data Fig. 3b), typically for high similarity segment duplications.

The current Shasta implementation does not attempt to remove the obviously incorrect connections. This results in unnecessary breaks in assembly contiguity. Despite this, Shasta assembly contiguity is adequate and comparable to what other, less performant long-read assemblers achieve. It is hoped that future Shasta releases will do a better job at handling these situations.

Marker graph. Consider a read whose marker representation is as follows:

a b c d e

We can represent this read as a directed graph that describes the sequence in which its markers appear. This is not very useful but illustrates the simplest form of a marker graph as used in the Shasta assembler. The marker graph is a directed graph in which each vertex represents a marker and each edge represents the transition between consecutive markers. We can associate sequence with each vertex and edge of the marker graph:

- Each vertex is associated with the sequence of the corresponding marker.
- If the markers of the source and target vertex of an edge do not overlap, the edge is associated with the sequence intervening between the two markers.
- If the markers of the source and target vertex of an edge do overlap, the edge is associated with the overlapping portion of the marker sequences.

Consider now a second read with the following marker representation, which differs from the previous one just by replacing marker c with x:

a b x d e

The marker graph for the two reads is Extended Data Fig. 4a. In the optimal alignment of the two reads, markers a, b, d and e are aligned. We can redraw the marker graph grouping together vertices that correspond to aligned markers as in

Extended Data Fig. 4b. Finally, we can merge aligned vertices to obtain a marker graph describing the two aligned reads, shown in Extended Data Fig. 4c.

Here, by construction, each vertex still has a unique sequence associated with it: the common sequence of the markers that were merged (however, the corresponding repeat counts can be different for each contributing read). An edge, on the other hand, can have different sequences associated with it; one corresponding to each of the contributing reads. In this example, edges $a \rightarrow b$ and $d \rightarrow e$ have two contributing reads, which can each have a distinct sequence between the two markers. We call coverage of a vertex or edge the number of reads ‘contributing’ to it. In this example, vertices a, b, d and e have coverage 2 and vertices c and x have coverage 1. Edges $a \rightarrow b$ and $d \rightarrow e$ have coverage 2, and the remaining edges have coverage 1.

The construction of the marker graph was illustrated here for two reads, but the Shasta assembler constructs a global marker graph that takes into account all oriented reads:

- (1) The process starts with a distinct vertex for each marker of each oriented read. Note that at this stage the marker graph is large ($\sim 2 \times 10^{10}$ vertices for a human assembly using default assembly parameters).
- (2) For each marker alignment corresponding to an edge of the read graph, we merge vertices corresponding to aligned markers.
- (3) Of the resulting merged vertices, we remove those whose coverage is too low or too high, indicating that the contributing reads or some of the alignments involved are probably in error. This is controlled by assembly parameters `MarkerGraph.minCoverage` (default 10) and `MarkerGraph.maxCoverage` (default 100), which specify the minimum and maximum coverage for a vertex to be kept.
- (4) Edges are created. An edge $v_0 \rightarrow v_1$ is created if there is at least a read contributing to both v_0 and v_1 and for which all markers intervening between v_0 and v_1 belong to vertices that were removed.

Note that this does not mean that all vertices with the same marker sequence are merged: two vertices are only merged if they have the same marker sequence, and if there are at least two reads for which the corresponding markers are aligned.

Given the large number of initial vertices involved, this computation is not trivial. To allow efficient computation in parallel on many threads a lock-free implementation of the disjoint data set data structure³³, is used for merging vertices. Some code changes were necessary to permit large numbers of vertices, as the initial implementation by Wenzel Jakob only allowed for 32-bit vertex identities (<https://github.com/wjakob/dset>).

Assembly graph. The Shasta assembly process also uses a compact representation of the marker graph, called the assembly graph, in which each linear sequence of edges is replaced by a single edge (Extended Data Fig. 5).

The length of an edge of the assembly graph is defined as the number of marker graph edges that it corresponds to. For each edge of the assembly graph, an average coverage is also computed, by averaging the coverage of the marker graph edges it corresponds to.

Using the marker graph to assemble sequence. The marker graph is a partial description of the multiple sequence alignment between reads and can be used to assemble consensus sequence. One simple way to do that is to only keep the ‘dominant’ path in the graph, and then traverse that path from vertex to edge to vertex, assembling a run-length encoded sequence as follows:

- (1) On a vertex, all reads have the same sequence, by construction: the marker sequence associated with the vertex. There is trivial consensus among all the reads contributing to a vertex, and the marker sequence can be used directly as the contribution of the vertex to assembled sequence.
- (2) For edges, there are two possible situations plus a hybrid case:
 - (a) If the adjacent markers overlap, in most cases all contributing reads have the same number of overlapping bases between the two markers, and we are again in a situation of trivial consensus, where all reads contribute the same sequence, which also agrees with the sequence of adjacent vertices. In cases where not all reads are in agreement on the number of overlapping bases, only reads with the most frequent number of overlapping bases are taken into account.
 - (b) If the adjacent markers do not overlap, then each read can have a different sequence between the two markers. In this situation, we compute a multiple sequence alignment of the sequences and a consensus using the `spoa` library⁴⁰ (<https://github.com/rvaser/spoa>). The multiple sequence alignment is computed constrained at both ends, because all reads contributing to the edge have, by construction, identical markers at both sides.
 - (c) A hybrid situation occasionally arises, in which some reads have the two markers overlapping, and some do not. In this case we count reads of the two kinds and discard the reads of the minority kind, then revert to one of the two cases 2(a) or 2(b) above.

This is the process used for sequence assembly by the current Shasta implementation. It requires a process to select and define dominant paths, which is described in the section ‘Selecting assembly paths in Shasta’. It is algorithmically

simple, but its main shortcoming is that it does not use for assembly reads that contribute to the abundant side branches. This means that coverage is lost, and therefore the accuracy of assembled sequence is not as good as it could be if all available coverage was used. Means to eliminate this shortcoming and use information from the side branches of the marker graph could be a subject of future work on the Shasta assembler.

This process described works with a run-length encoded sequence and therefore assembles a run-length encoded sequence. The final step to create raw assembled sequence is to compute the most likely repeat count for each sequence position in RLD. After some experimentation, this is currently done by choosing as the most likely repeat count the one that appears the most frequently in the reads that contributed to each assembled position.

A simple Bayesian model for repeat counts resulted in a modest improvement in the quality of assembled sequence. But the model appears to sensitive to calibration errors, and therefore it is not used by default in Shasta assemblies. However, it is used by `MarginPolish`, as described in the `MarginPolish` section.

Selecting assembly paths in Shasta. The sequence assembly procedure described in the previous section can be used to assemble sequence for any path in the marker graph. This section describes the selection of paths for assembly in the current Shasta implementation. This is done by a series of steps that ‘remove’ edges (but not vertices) from the marker graph until the marker graph consists mainly of linear sections that can be used as the assembly paths. For speed, edges are not actually removed but just marked as removed using a set of flag bits allocated for this purpose in each edge. However, the description that follows will use the loose term ‘remove’ to indicate that an edge was flagged as removed.

This process consists of the following three steps, described in more detail in the following sections:

- (1) Approximate transitive reduction of the marker graph
- (2) Pruning of short side branches (leaves)
- (3) Removal of bubbles and superbubbles

The last step, removal of bubbles and superbubbles, is consistent with Shasta’s current assembly goal, which is to compute a mostly monoploid assembly, at least on short scales.

Approximate transitive reduction of the marker graph. The goal of this step is to eliminate the side branches in the marker graph, which are the result of errors. Despite the fact that the number of side branches is substantially reduced thanks to the use of RLE, side branches are still abundant. This step uses an approximate transitive reduction of the marker graph that only considers reachability up to a maximum distance, controlled by assembly parameter `MarkerGraph.maxDistance` (default 30 marker graph edges). Using a maximum distance makes sure that the process remains computationally affordable, and also has the advantage of not removing long-range edges in the marker graph, which could be substantial.

In detail, the process works as follows. In this description, the edge being considered for removal is the edge $v_0 v_1$ with source vertex v_0 and target vertex v_1 . The first two steps are not really part of the transitive reduction but are performed by the same code for convenience.

- (1) All edges with coverage less than or equal to `MarkerGraph.lowCoverageThreshold` are unconditionally removed. The default value for this assembly parameter is 0, so this step does nothing when using default parameters.
- (2) All edges with coverage 1 and for which the only supporting read has a large marker skip are unconditionally removed. The marker skips of an edge, for a given read, is defined as the distance (in markers) between the v_0 marker for that read and the v_1 marker for the same read. Most marker skips are small, and a large skip is indicative of an artifact. Keeping those edges could result in assembly errors. The marker skip threshold is controlled by assembly parameter `MarkerGraph.edgeMarkerSkipThreshold` (default 100 markers).
- (3) Edges with coverage greater than `MarkerGraph.lowCoverageThreshold` (default 0) and less than `MarkerGraph.highCoverageThreshold` (default 256), and that were not previously removed, are processed in order of increasing coverage. Note that with the default values of these parameters all edges are processed, because edge coverage is stored using 1 byte and therefore can never be more than 255 (it is saturated at 255). For each edge $v_0 v_1$, a breadth-first search (BFS) in the alternative path from v_0 to v_1 exists, edge $v_0 v_1$ is removed. Note that the BFS does not use edges that have already been removed, and so the process is guaranteed not to affect reachability. Processing edges in order of increasing coverage makes sure that low coverage edges the most likely to be removed.

The transitive reduction step is intrinsically sequential and so it is currently performed in sequential code for simplicity. It could be parallelized in principle, but that would require sophisticated locking of marker graph edges to make sure independent threads do not step on each other, possibly reducing reachability. However, even with sequential code, this step is not computationally expensive, taking typically only a small fraction of total assembly time.

When the transitive reduction step is complete, the marker graph consists mostly of linear sections composed of vertices with an in-degree and out-degree of

one, with occasional side branches and bubbles or superbubbles, which are handled in the next two phases described in the following.

Pruning of short side branches (leaves). At this stage, a few iterations of pruning are done by simply removing, at each iteration, edge $v_0 v_1$ if v_0 has in-degree 0 (that is, is a backward-pointing leaf) or v_1 has out-degree 0 (that is, is a forward-pointing leaf). The net effect is that all side branches of length (number of edges) at most equal to the number of iterations are removed. This leaves the leaf vertex isolated, which causes no problems. The number of iterations is controlled by assembly parameter `MarkerGraph.pruneIterationCount` (default 6).

Removal of bubbles and superbubbles. The marker graph now consists of mostly linear section with occasional bubbles or superbubbles⁶⁴. Most of the bubbles and superbubbles are caused by errors, but some of those are due to heterozygous loci in the genome being assembled. Bubbles and superbubbles of the latter type could be used for separating haplotypes (phasing), a possibility that will be addressed in future Shasta releases. However, the goal of the current Shasta implementation is to create a monoploid assembly at all scales but the very long ones. Accordingly, bubbles and superbubbles at short scales are treated as errors, and the goal of the bubble/superbubble removal step is to keep the most significant path in each bubble or superbubble. The Extended Data Fig. 6 shows typical examples of a bubble and superbubble in the marker graph.

The bubble/superbubble removal process is iterative. Early iterations work on short scales, and late iterations work on longer scales. Each iteration uses a length threshold that controls the maximum number of marker graph edges for features to be considered for removal. The value of the threshold for each iteration is specified using assembly parameter `MarkerGraph.simplifyMaxLength`, which consists of a comma-separated string of integer numbers, each specifying the threshold for one iteration in the process. The default values are 10, 100 and 1,000, which means that three iterations of this process are performed. The first iteration uses a threshold of ten marker graph edges, and the second and third iterations use length thresholds of 100 and 1,000 marker graph edges, respectively. The last and largest of the threshold values used determines the size of the smallest bubble or superbubble that will survive the process. The default 1,000 markers are equivalent to roughly 13 kb. To suppress more bubble/superbubbles, increase the threshold for the last iteration. To see more bubbles/superbubbles, decrease the length threshold for the last iteration or remove the last iteration entirely.

The goal of the increasing threshold values is to work on small features at first, and on larger features in the later iterations. The choice of `MarkerGraph.simplifyMaxLength` could be application dependent. The default value is a reasonable compromise useful if one desires a mostly monoploid assembly with just some large heterozygous features.

Each iteration consists of two steps. The first removes bubbles and the second removes superbubbles. Only bubbles/superbubbles consisting of features shorter than the threshold for the current iteration are considered:

- (1) Bubble removal
 - (a) An assembly graph corresponding to the current marker graph is created.
 - (b) Bubbles are located in which the length of all branches (number of marker graph edges) is no more than the length threshold at the current iteration. In the assembly graph, a bubble appears as a set of parallel edges (edges with the same source and target).
 - (c) In each bubble, only the assembly graph edge with the highest average coverage is kept. Marker graph edges corresponding to all other assembly graph edges in the bubble are flagged as removed.
- (2) Superbubble removal
 - (a) An assembly graph corresponding to the current marker graph is created.
 - (b) Connected components of the assembly graph are computed, but only considering edges below the current length threshold. This way, each connected component corresponds to a 'cluster' of 'short' assembly graph edges.
 - (c) For each cluster, entries in the cluster are located. These are vertices that have in-edges from a vertex outside the cluster. Similarly, out-edges are located (vertices that have out-edges outside the cluster).
 - (d) For each entry/exit pair, the shortest path is computed. However, in this case the 'length' of an assembly graph edge is defined as the inverse of its average coverage: that is, the inverse of average coverage for all the contributing marker graph edges.
 - (e) Edges on each shortest path are marked as edges to be kept.
 - (f) All other edges internal to the cluster are removed.

When all iterations of bubble/superbubble removal are complete, the assembler creates a final version of the assembly graph. Each edge of the assembly graph corresponds to a path in the marker graph, for which sequence can be assembled using the method described. Note, however, that the marker graph and the assembly graph have been constructed to contain both strands. Special care is taken during all transformation steps to make sure that the marker graph (and therefore

the assembly graph) remain symmetric with respect to strand swaps. Therefore, most assembly graph edges come in reverse complemented pairs, of which we assemble only one. However, it is possible but rare for an assembly graph to be its own reverse complement.

Assembly parameters selection. The sequence of computational steps outlined before depends on a number of assembly parameters, such as, for example, the length and fraction of k -mers used as markers, the parameters controlling the LowHash iteration and so on. In Shasta, all of these parameters are exposed as command line options and none of them are hardcoded or hidden. Our error analysis shows that the set of assembly parameters we used (the default values for Shasta v.0.1.0) gave satisfactory assembly results for our data. However, we do not claim that the same choices would generalize to other situations. Additional work will be needed to find parameter sets that work for lower or higher coverage, for genomes of different sizes and characteristics or for different types of long read.

High performance computing techniques used by Shasta. The Shasta assembler is designed to run on a single machine with an amount of memory sufficient to hold all of its data structures (1–2 Tb for a human assembly, depending on coverage). All data structures are memory mapped and can be set up to remain available after assembly completes. Note that using such a large memory machine does not substantially increase the cost per CPU cycle. For example, on AWS the cost per virtual processor hour for large memory instances is no more than twice the cost for laptop-sized instances.

There are various advantages to running assemblies in this way:

- Running on a single machine simplifies the logistics of running an assembly versus, for example, running on a cluster of smaller machines with shared storage.
- No disk input/output takes place during assembly, except for loading the reads in memory and writing out assembly results plus a few small files containing summary information. This eliminates performance bottlenecks commonly caused by disk I/O.
- Having all data structures in memory makes it easier and more efficient to exploit parallelism, even at very low granularity.
- Algorithm development is easier, as all data are immediately accessible without the need to read files from disk. For example, it is possible to easily rerun a specific portion of an assembly for experimentation and debugging without any wait time for data structures to be read from disk.
- When the assembler data structures are set up to remain in memory after the assembler completes, it is possible to use the Python API or the Shasta http server to inspect and analyze an assembly and its data structures (for example, display a portion of the read graph, marker graph or assembly graph).
- For optimal performance, assembler data structures can be mapped to Linux 2 MB pages ('huge pages'). This makes it faster for the operating system to allocate and manage the memory, and improves translation lookaside buffer efficiency. Using huge pages mapped on the `hugetlbfs` filesystem (Shasta executable options `--memoryMode filesystem --memoryBacking 2M`) can result in a notable speedup (20–30%) for large assemblies. However, it requires root privilege via `sudo`.

To optimize performance in this setting, the Shasta assembler uses various techniques:

- In most parallel steps, the division of work among threads is not set up in advance but decided dynamically ('dynamic load balancing'). As a thread finishes a piece of work assigned to it, it grabs another chunk of work to do. The process of assigning work items to threads is lock-free (that is, it uses atomic memory primitives rather than mutexes or other synchronization methods provided by the operating system).
- Most large memory allocations are done via `mmap` and can optionally be mapped to Linux 2 MB pages backed by the Linux `hugetlbfs`. This memory is persistent until the next reboot and is resident (nonpageable). As a result, assembler data structures can be kept in memory and reaccessed repeatedly at very low cost. This facilitates algorithm development (for example, it allows repeatedly testing a single assembly phase without having to rerun the entire assembly each time or having to wait for data to load) and postprocessing (inspecting assembly data structures after the assembly is complete). The Shasta http server and Python API take advantage of this capability.
- The Shasta code includes a C++ class for conveniently handling these large memory-mapped regions as C++ containers with familiar semantics (`class shasta::MemoryMapped::Vector`).
- In situations where a large number of small vectors are required, a two-pass process is used (`class shasta::MemoryMapped::VectorOfVectors`). In the first pass, one computes the length of each of the vectors. A single large area is then allocated to hold all of the vectors contiguously, together with another area to hold indexes pointing to the beginning of each of the short vectors. In a second pass, the vectors are then filled. Both passes can be performed in parallel and are entirely lock free. This process eliminates memory allocation overhead that would be incurred if each of the vectors were to be allocated individually.

Thanks to these techniques, Shasta achieves close to 100% CPU use during its parallel phases, even when using large numbers of threads. However, a number of sequential phases remain, which typically result in average CPU use during a large assembly around 70%. Some of these sequential phases can be parallelized, which would result in increased average CPU use and improved assembly performance.

MarginPolish. Throughout, we used MarginPolish v.1.0.0 from <https://github.com/ucsc-nanopore-cgl/MarginPolish>.

MarginPolish is an assembly refinement tool designed to sum over (marginalize) read to assembly alignment uncertainty. It takes as input a genome assembly and set of aligned reads in BAM format. It outputs a refined version of the input genome assembly after attempting to correct base-level errors in terms of substitutions and indels (insertions and deletions). It can also output a summary representation of the assembly and read alignments as a weighted POA graph, which is used by the HELEN neural network-based polisher described next.

It was designed and is optimized to work with noisy long ONT reads, although parameterization for other, similar read types is easily possible. It does not yet consider signal-level information from ONT reads. It is also currently a haploid polisher; in that it does not attempt to recognize or represent heterozygous polymorphisms or phasing relationships. For haploid genome assemblies of a diploid genome, it will therefore fail to capture half of all heterozygous polymorphisms.

Algorithm overview. MarginPolish works as follows:

- Reads and the input assembly are converted to their RLE (see Shasta the description for the steps and rationale).
- A restricted, weighted POA⁴⁰ graph is constructed representing the RLE input assembly and potential edits to it in terms of substitutions and indels.
- Within identified regions of the POA containing likely assembly errors:
 - A set of alternative sequences representing combinations of edits are enumerated by locally traversing the POA within the region.
 - The likelihood of the existing and each alternative sequence is evaluated given the aligned reads.
 - If an alternative sequence with higher likelihood than the current reference exists, then the assembly at the location is updated with this higher likelihood sequence.
- Optionally, the program loops back to step 2 to repeat the refinement process (by default it loops back once).
- The modified run-length encoded assembly is expanded by estimating the repeat count of each base given the reads using a simple Bayesian model. The resulting final, polished assembly is output. In addition, a representation of the weighted POA can be output.

Innovations. Compared to existing tools, MarginPolish is most similar to Racon⁴² in that they are comparable in speed, both principally use small-parameter HMM-like models and both do not currently use signal information. Compared to Racon, MarginPolish has some key innovations that we have found to improve polishing accuracy:

- MarginPolish, as with our earlier tool in the Margin series², uses the forward-backward and forward algorithms for pair-HMMs to sum over all possible pairwise alignments between pairs of sequences instead of the single most probable alignment (Viterbi). Considering all alignments allows more information to be extracted per read.
- The POA graph is constructed from a set of weights computed from the posterior alignment probabilities of each read to the initial assembled reference sequence (see below), the result is that MarginPolish POA construction does not have a read-order dependence. This is similar to that described by HGAP3 (ref. ⁶⁵). Most earlier algorithms for constructing POA graphs have a well-known explicit read-order dependence that can result in undesirable topologies⁴⁰.
- MarginPolish works in a run-length encoded space, which results in considerably less alignment uncertainty and correspondingly improved performance.
- MarginPolish, similarly to Nanopolish⁴⁶, evaluates the likelihood of each alternative sequence introduced into the assembly. This improves performance relative to a faster but less accurate algorithm that traces back a consensus sequence through the POA graph.
- MarginPolish uses a simple chunking scheme to break up the polishing of the assembly into overlapping pieces. This results in low memory usage per core and simple parallelism.

In the following, steps 2, 3 and 5 of the MarginPolish algorithm are described in detail. In addition, the parallelization scheme is described.

POA graph construction. To create the POA, we start with the existing assembled sequence $s = s_1, s_2, \dots, s_n$, and for each read $r = r_1, r_2, \dots, r_m$ in the set of reads R use the forward-backward algorithm with a standard three-state, affine-gap pair-HMM to derive posterior alignment probabilities using the implementation described in ref. ⁵⁶. The parameters for this model are specified in the `polish.hmm` subtree of

the JSON formatted parameters file, including `polish.hmm.transitions` and `polish.hmm.emissions`. Current defaults were tuned via expectation maximization¹² of R9.4 ONT reads aligned to a bacterial reference; we have observed that the parameters for this HMM seem robust to small changes in basecaller versions. The result of running the forward-backward algorithm is three sets of posterior probabilities:

- First, match probabilities: the set of posterior match probabilities, each the probability $P(r_i \triangleleft s_j)$ that a read base r_i is aligned to a base s_j in s .
- Second, insertion probabilities: the set of posterior insertion probabilities, each the probability $P(r_i \triangleleft -j)$ that a read base r_i is inserted between two bases s_j and s_{j+1} in s , or, if $j=0$, inserted before the start of s , or, if $j=n$, after the end of s .
- Third, deletion probabilities, the set of posterior deletion probabilities, each the probability $P(-i \triangleleft s_j)$ that a base s_j in s is deleted between two read bases r_i and r_{i+1} . (Note that because a read is generally an incomplete observation of s , we consider the probability that a base in s is deleted before the first position or after the last position of a read as 0.)

As most probabilities in these three sets are very small and yet to store and compute all the probabilities would require evaluating comparatively large forward and backward alignment matrices, we restrict the set of probabilities heuristically as follows:

- We use a banded forward-backward algorithm, as originally described in ref. ⁶⁷. To do this we use the original alignment of the read to s as in the input BAM file. Given that s is generally much longer than each read this allows computation of each forward-backward invocation in time linearly proportional to the length of each read, at the cost of restricting the probability computation to a subportion of the overall matrix, albeit one that contains most of the probability mass.
- We only store posterior probabilities above a threshold (`polish.pairwiseAlignmentParameters.threshold`, by default 0.01), treating smaller probabilities as equivalent as zero.

The result is that these three sets of probabilities are a very sparse subset of the complete sets.

To estimate the posterior probability of a multi-base insertion of a read substring r_p, r_{p+1}, \dots, r_k at a given location j in s involves repeated summation over terms in the forward and backward matrices. Instead, to approximate this probability we heuristically use

$$P(r_i, r_{i+1}, \dots, r_k \diamond -j) = \underset{l \in [i, k]}{\operatorname{argmin}} P(\eta \diamond -j)$$

the minimum probability of any base in the multi-base insertion being individually inserted at the location in s as a proxy, a probability that is an upper bound on the actual probability.

Similarly, we estimate the posterior probability of a deletion involving more than one contiguous base s at a given location in a read using analogous logic. As we store a sparse subset of the single-base insertion and deletion probabilities and given these probability approximations, it is easy to calculate all the multi-base indel probabilities with value greater than t by linear traversal of the single-based insertion and deletion probabilities after sorting them, respectively, by their read and s coordinates. The result of such calculation is expanded sets of insertion and deletion probabilities that include multi-base probabilities.

To build the POA we start from s , which we call the backbone. The backbone is a graph where each base s_j in s corresponds to a node, there are special source and sink nodes (which do not have a base label), and the directed edges connect the nodes for successive bases s_j, s_{j+1} in s , from the source node to the node for s_j , and, similarly, from the node for s_j to the sink node.

Each nonsource/sink node in the backbone has a separate weight for each possible base $x \in \{A, C, G, T\}$. This weight (w) is

$$w(j, x) = \sum_{r \in R} \sum_i 1_x(r_i) P(r_i \diamond s_j)$$

where $1_x(r_i)$ is an indicator function that is 1 if $r_i = x$ and otherwise 0, corresponds to the sum of match probabilities of read elements of base x being aligned to s_j . This weight has a probabilistic interpretation: it is the total number of expected observations of the base x in the reads aligned to s_j , summing over all possible pairwise alignments of the reads to s . It can be fractional because of the inherent uncertainty of these alignments; for example, we may predict only a 50% probability of observing such a base in a read.

We add deletion edges that connect nodes in the backbone. Indexing the nodes in the backbone from 0 (the source) to the source $n+1$ (the sink), a deletion edge between positions j and k in the backbone corresponds to the deletion of bases $j, j+1, \dots, k$ in s . Each deletion edge has a weight equal to the sum of deletion probabilities for deletion events that delete the corresponding base(s) in s , summing over all possible deletion locations in all reads. Deletions with no weight are not included. Again, this weight has a probabilistic interpretation: it is the expected number of times we see the deletion in the reads, and again it may be fractional.

We represent insertions as nodes labeled with an insertion sequence. Each insertion node has a single incoming edge from a backbone node, and a single outgoing edge to the next backbone node in the backbone sequence. Each insertion is labeled with a weight equal to the sum of probabilities of events that insert the given insertion sequence between the corresponding bases in s . The resulting POA is a restricted form of a weighted, directed acyclic graph (Extended Data Fig. 7a shows an example)

Frequently, either an insertion or deletion can be made between different successive bases in s resulting in the same edited sequence. To ensure that such equivalent events are not represented multiple times in the POA, and to ensure we sum their weights correctly, we 'left shift' indels to their maximum extent. When shifting an indel results in multiple equivalent deletion edges or insertions, we remove the duplicate elements, updating the weight of the residual element to include the sum of the weights of the removed elements. For example, the insertion of 'AT' in Extended Data Fig. 7 is shifted left to its maximal extent and could include the merger of an equivalent 'AT' insertion starting two backbone nodes to the right.

Local haplotype proposal. After constructing the POA we use it to sample alternative assemblies. We first prune the POA to mark indels and base substitutions with weight below a threshold, which are generally the result of sequencing errors (Extended Data Fig. 7b). Currently, this threshold (polish.candidateVariantWeight=0.18, established empirically) is normalized as a fraction of the estimated coverage at the site, which is calculated in a running window around each node in the backbone of 100 bases. Consequently, if fewer than 18% of the reads are expected to include the change then the edit is pruned from consideration.

To further avoid a combinatorial explosion, we sample alternative assemblies locally. We identify subgraphs of s containing indels and substitutions to s then in each subgraph, defined by a start and end backbone vertex, we enumerate all possible paths between the start and end vertex and all plausible base substitutions from the backbone sequence. The rationale for heuristically doing this locally is that two subgraphs separated by one or more anchor backbone sites with no plausible edits are conditionally independent of each other given the corresponding interstitial anchoring substring of s and the substrings of the reads aligning to it. Currently, any backbone site more than polish.columnAnchorTrim=5 nodes (equivalent to bases) in the backbone from a node overlapping a plausible edit (either substitution or indel) is considered an anchor. This heuristic allows for some exploration of alignment uncertainty around a potential edit. Given the set of anchors computation proceeds by identifying successive pairs of anchors separated by subgraphs containing the potential edits, with the two anchors considered the source and sink vertex.

A simple Bayesian model for run-length decoding. RLE allows for separate modeling of length and nucleotide error profiles. In particular, length predictions are notoriously error prone in nanopore basecalling. Since homopolymers produce continuous signals, and DNA translocates at a variable rate through the pore, the basecaller often fails to infer the true number of bases given a single sample. For this reason, a Bayesian model is used for error correction in the length domain, given a distribution of repeated samples at a locus.

To model the error profile, a suitable reference sequence is selected as the truth set. Reads and reference are run-length encoded and aligned by their nucleotides. The alignment is used to generate a mapping of observed lengths to their true length (y, x) where y = true and x = observed for each position in the alignment. Observations from alignment are tracked using a matrix of predefined size ($y_{\max} = 50, x_{\max} = 50$) in which each coordinate contains the corresponding count for (y, x). Finally, the matrix is normalized along one axis to generate a probability distribution of $P(X|y_j)$ for j in $[1, y_{\max}]$. This process is performed for each of the four bases.

With enough observations, the model can be used to find the most probable true run length given a vector of observed lengths X . This is done using a simple log likelihood calculation over the observations x_i for all possible true lengths y_j in Y , assuming the length observations to be independent and identically distributed. The length y_j corresponding to the greatest likelihood $P(X|y_j, \text{Base})$ is chosen as the consensus length for each alignment position (Extended Data Fig. 8).

Training. To generate a model, we ran MarginPolish with reads from a specific basecaller version aligned to a reference (GRCh38) and specified the -outputRepeatCounts flag. This option produces a tab-separated value for each chunk describing all the observed repeat counts aligned to each backbone node in the POA. These files are consumed by a script in the https://github.com/rlorigro/runlength_analysis repository, which generates a run-length encoded consensus sequence, aligns to the reference and performs the described process to produce the model.

The allParams.np.human.guppy-ff-235.json model used for most of the analysis was generated from HG00733 reads basecalled with Guppy Flipflop v.2.3.5 aligned to GRCh38, with chromosomes 1, 2, 3, 4, 5, 6 and 12 selected. The model allParams.np.human.guppy-ff-233.json was generated from Guppy Flipflop v.2.3.3 data and chromosomes 1–10 were used. This model was also used for the CHM13 analysis, as the run-length error profile is very similar between v.2.3.3 and v.2.3.1.

Parallelization and computational considerations. To parallelize MarginPolish we break the assembly up into chunks of size polish.chunkSize=1000 bases, with an overlap of polish.chunkBoundary=50 bases. We then run the MarginPolish algorithm on each chunk independently and in parallel, stitching together the resulting chunks after finding an optimal pairwise alignment (using the default HMM described earlier) of the overlaps that we use to remove the duplication.

Memory usage scales with thread count, read depth, and chunk size. For this reason, we downsample reads in a chunk to polish.maxDepth=50 \times coverage by counting total nucleotides in the chunk N_c and discarding reads with likelihood $1 - (\text{chunkSize} + 2 \times \text{chunkBoundary}) \times \text{maxDepth} / N_c$. With these parameters, we find that 2 GB of memory per thread is sufficient to run MarginPolish on genome-scale assemblies. Across 13 whole-genome runs, we averaged roughly 350 CPU hours per Gb of assembled sequence.

HELEN. HELEN is a deep neural network-based haploid consensus sequence polisher. HELEN uses a multi-task RNN⁴¹ that takes the weights of the POA graph of MarginPolish to predict a base and a run length for each genomic position. MarginPolish constructs the POA graph by performing multiple possible alignments of a single read that makes the weights associative to the correct underlying base and a run length. The RNN used in HELEN takes advantage of the transitive relationship of the genomic sequence and associative coupling of the POA weights to the correct base and run length to produce a consensus sequence with higher accuracy.

The error correction with HELEN is done in three steps. First, we generate tensor-like images of genomic segments with MarginPolish that encodes POA graph weights for each genomic position. Then we use a trained RNN model to produce predicted bases and run lengths for each of the generated images. Finally, we stitch the chunked sequences to get a contiguous polished sequence.

Image generation. MarginPolish produces an image-like summary of the final POA state for use by HELEN. At a high level, the image summarizes the weighted alignment likelihoods of all reads divided into nucleotide, orientation and run length.

The positions of the POA nodes are recorded using three coordinates: the position in the backbone sequence of the POA, the position in the insert sequences between backbone nodes and the index of the run-length block. All backbone positions have an insert coordinate of 0. Each backbone and insert coordinate include one or more run-length coordinate.

When encoding a run length we divide all read observations into blocks from 0 to 10 inclusive (this length is configurable). For cases where no observations exceed the maximum run length, a single run-length image can describe the POA node. When an observed run length exceeds the length of the block, the run length is encoded as that block's maximum (10) and the remaining run length is encoded in successive blocks. For a run length that terminates in a block, its weight is contributed to the run-length 0 column in all successive blocks. This means that the records for all run-length blocks of a given backbone and insert position have the same total weight. As an example, consider three read positions aligned to a node with run lengths of 8, 10 and 12. These require two run-length blocks to describe: the first block includes one 8 and two 10s, and the second includes two 0s and one 2.

The information described at each position (backbone, insert and run length) is encoded in 92 features: each nucleotide {A, C, T, G} and run length {0, 1, ..., 10}, plus a gap weight (for deletions in read alignments). The weights for each of these 45 observations are separated into forward and reverse strand for a total of 90 features. The weights for each of these features are normalized over the total weight for the record and accompanied by an additional data point describing the total weight of the record. This normalization column for the record is an approximation of the read depth aligned to that node. Insert nodes are annotated with a binary feature (for a final total of 92); weights for an insert node's alignments are normalized over total weight at the backbone node it is rooted at (not the weight of the insert node itself) and gap alignment weights are not applied to them.

Labeling nodes for training require a truth sequence aligned to the assembly reference. This provides a genome-scale location for the true sequence and allows its length to help in the resolution of segmental duplications or repetitive regions. When a region of the assembly is analyzed with MarginPolish, the truth sequences aligned to that region are extracted. If there is not a single truth sequence that roughly matches the length of the consensus for this region, we treat it as an uncertain region and no training images are produced. Having identified a suitable truth sequence, it is aligned to the final consensus sequence in nonrun-length space with Smith-Waterman. Both sequences and the alignment are then run-length encoded, and true labels are matched with locations in the images. All data between the first and last matched nodes are used in the final training images (leading and trailing inserts or deletes are discarded). For our training, we aligned the truth sequences with minimap2 using the asm20 preset and filtered the alignments to include only primary and supplementary alignments (no secondary alignments).

Extended Data Fig. 9 shows a graphical representation of the images. On the y axis, we display true nucleotide labels (with the dash representing no alignment or gap) and true run length. On the x axis, the features used as input to HELEN are displayed: first, the normalization column (the total weight at the backbone

position); second, the insert column (the binary feature encoding whether the image is for a backbone or insert node); 48 columns describing the weights associated with read observations (stratified by nucleotide, run length, strand) and two columns describing weights for gaps in read alignments (stratified by strand). In this example, we have reduced the maximum run length per block from ten to five for demonstration purposes.

We selected these two images to highlight three features of the model: the way multiple run-length blocks are used to encode observations for a single node, and the relevant features around a true gap and a true insert that enable HELEN to correct these errors.

To illustrate multiple run-length blocks, we highlight two locations on the image Extended Data Fig. 9a(i). The first are the nodes labeled (A,5) and (A,3). This is the labeling for a true (A,8) sequence separated into two blocks. See that the bulk of the weight is on the (A,5) features on the first block, with most of that distributed across the (A,1–3) features on the second. Second, observe the nodes on Extended Data Fig. 9a(i) labeled (T,4) and (T,0). Here, we show the true labeling of a (T,4) sequence where there are some read observations extending into a second run-length block.

To show a feature of a true gap, note on Extended Data Fig. 9a(i) the noninsert nodes labeled (–,0). We know that MarginPolish predicted a single cytosine nucleotide (as it is a backbone node and the (C,1) nodes have the bulk of the weight. Here, HELEN is able to use the low overall weight (the lighter region in the normalization column) at this location as evidence of fewer supporting read alignments and can correct the call.

The position labeled (G,2) on Extended Data Fig. 9a(i) details a true insertion. It is not detected by MarginPolish (as all insert nodes are not included in the final consensus sequence). Read support is present for the insert, less than the backbone nodes in this image but more than the other insert nodes. HELEN can identify this sequence and correct it.

Finally, we note that the length of the run-length blocks results in streaks at multiples of this length (10) for long homopolymers. The root of this effect lies in the basecaller producing similar prediction distributions for these cases (that is, the run-length predictions made by the basecaller for a true run length of 25 are similar to the run-length predictions made for a true run length of 35, see Fig. 4b Guppy v.2.3.3). This gives the model little information to differentiate on, and the issue is exacerbated by the low occurrence of long run lengths in the training data. Because the model divides run-length observations into chunks of size 10, it tends to call the first chunks correctly (having length 10) but has very low signal for the last chunk and most often predicts 0.

The model. We use a sequence transduction model for consensus polishing. The model structure consists of two single layer gated recurrent neural units (GRU) for encoding and decoding on top of two linear transformation layers (Extended Data Fig. 10). The two linear transformation layers independently predict a base and a run length for each position in the input sequence. Each unit of the GRU can be described using the four functions it calculates:

$$\begin{aligned} r_t &= \text{Sigmoid}(W_{ir}x_t + W_{ir}h_{t-1}) \\ u_t &= \text{Sigmoid}(W_{iu}x_t + W_{iu}h_{t-1}) \\ n_t &= \tanh(W_{in}x_t + r_t * (W_{in}h_{t-1})) \\ h_t &= (1 - u_t) * n_t + u_t * h_{t-1} \end{aligned}$$

For each genomic position t , we calculate the current state h_t from the new state n_t and the update value u_t applied to the output state of previous genomic position h_{t-1} . The update function u_t decides how much past information to propagate to the next genomic position. It multiplies the input x_t with the weight vector W_{iu} and multiplies the hidden state of the previous genomic position h_{t-1} . The weight vectors decide how much from the previous state to propagate to the next state. The reset function r_t decides how much information to dissolve from the previous state. Using a different weight vector, the r_t function decides how much information to dissolve from the past. The new memory state n_t is calculated by multiplying the input x_t with the weight vector W_{in} and applying a Hadamard multiplication $*$ between the reset function value and a weighted state of the previous hidden state h_{t-1} . The new state captures the associative relationship between the input weights and true prediction. In this setup, we can see that r_t and u_t can decide to hold memory from distant locations while n_t captures the associative nature of the weights to the prediction, helping the model to decide how to propagate genomic information in the sequence. The output of each genomic position h_t can be then fed to the next genomic position as a reference to the previously decoded genomic position. The final two layers apply linear transformation functions:

$$B_t = h_t * W^T$$

$$R_t = h_t * W^T$$

The two linear transformation functions independently calculate a base prediction B_t and a run-length prediction R_t from the hidden state output of that genomic position h_t . The model operates in hard parameter sharing mode where the model learns to perform two tasks using the same set of underlying

parameters from the GRU layers. The ability of the model to reduce the error rate of the assemblies from multiple samples with multiple assemblers shows the generalizability and robustness we achieve with this method.

Sliding window mechanism. One of the challenges of this setup is the sequence length. From the functions of recurrent units, we see that each state is updated based on the previous state and associated weight. Due to the noisy nature of the data, if the sequence length is too long, the back-propagation becomes difficult over noisy regions. On the other hand, a small sequence length would make the program very slow. We balance the runtime and accuracy by using a sliding window approach.

During the sliding window, we chunk the sequence of thousand bases to multiple overlapping windows of length 100. Starting from the leftmost window, we perform prediction on sequence pileups of the window and transmit the hidden state of the current window to the next window and slide the window by 50 bases to the right. For each window, we collect all the predicted values and add it to a global sequence inference counter that can keep track of predicted probabilities of base and run length at each position. Last, we aggregate the probabilities from the global inference counter to generate a sequence. This setup allows us to use the minibatch feature of the popular neural network libraries allowing inference on a batch of inputs instead of performing inference one at a time.

Training the model. HELEN is trained with a gradient descent method. We use an adaptive moment estimation (Adam) method to compute gradients for each of the parameters in the model based on a target loss function. Adam uses both decaying squared gradients and the decaying average of gradients, making it suitable to use with RNNs. Adam performs gradient optimization by adapting the parameters to set in a way that minimizes the value of the loss function.

We perform optimization through back-propagation per window of the input sequence. From the equations of the linear transformation function, we see that we get two vectors $\mathbf{B} = [B_1, B_2, B_3, \dots, B_n]$ and $\mathbf{R} = [R_1, R_2, R_3, \dots, R_n]$ containing base and run-length predictions for each window of size n . From the labeled data we get two more such vectors $\mathbf{T}_B = [T_{B1}, T_{B2}, T_{B3}, \dots, T_{Bn}]$ and \mathbf{T}_R containing the true base and true run-length encoded values of each position in the window. From these losses, the function of the loss L is calculated:

$$\begin{aligned} L_{B(\mathbf{B}, \mathbf{T}_B)} &= -\mathbf{B}[\mathbf{T}_B] + \log\left(\sum_j \exp(\mathbf{B}[j])\right) \\ L_{R(\mathbf{R}, \mathbf{T}_R)} &= \text{weight}[\mathbf{T}_R] \left(-\mathbf{R}[\mathbf{T}_R] + \log\left(\sum_j \exp(\mathbf{R}[j])\right) \right) \\ L &= L_B + L_R \end{aligned}$$

L_B calculates the base prediction loss and L_R calculates the run-length encoded prediction loss. The run-length encoded class distribution is heavily biased toward lower run-length values, so, we apply class-wise weights depending on the observation of per class to make the learning process balanced between classes. The optimizer then updates the parameters or weights of GRU layers and linear layers in a way that minimizes the value of the loss function. We can see that the loss function is a summation of the two independent loss functions but the underlying weights from the RNN belongs to the same set of elements in the model. In this setting, the model optimizes to learn both tasks simultaneously by updating the same set of weights.

Sequence stitching. To parallelize the polishing pipeline, MarginPolish chunks the genome into smaller segments while generating images. Each image segment encodes 1,000 nucleotide bases, and two adjacent chunks have 50 nucleotide bases of overlap between them. During the inference step, we save all run-length and base predictions of the images, including their start and end genomic positions.

For stitching, we load all the image predictions and sort them based on the genomic start position of the image chunk and stitch them in parallel processes. For example, if there are n predictions from n images of a contig and we have t available threads, we divide n prediction chunks into t buckets each containing approximately $\frac{n}{t}$ predicted sequences. Then we start t processes in parallel where each process stitches all the sequences assigned to it and returns a longer sequence. For stitching two adjacent sequences, we take the overlapping sequences between the two sequences and perform a pairwise Smith–Waterman alignment. From the alignment, we pick an anchor position where both sequences agree the most and create one sequence. After all the processes finish stitching the buckets, we get t longer sequences generated by each process. Finally, we iteratively stitch the t sequences using the same process and get one contiguous sequence for the contig.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

Sequence data including raw signal files (FAST5) and basecalls (FASTQ) are available as an AWS Open Data set for download from <https://github.com/human-pangenomics/hpgp-data>. Nanopore sequence data and polished assemblies are additionally archived and available from the European Nucleotide Archive

under accession code PRJEB37264. Source data for Figs. 1–5 are presented with the paper.

Code availability

Shasta (<https://github.com/chanzuckerberg/shasta>), MarginPolish (<https://github.com/UCSC-nanopore-cgl/marginPolish>) and HELEN (<https://github.com/kishwarshafin/helen>) are publicly available. They have open-source MIT licenses that fully support the open source initiative.

References

- Sedlazeck, F. J. et al. SVCollector: optimized sample selection for validating and long-read resequencing of structural variants. Preprint at *bioRxiv* <https://doi.org/10.1101/342386> (2018).
- Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- PacBio. Data release: highest-quality, most contiguous individual human genome assembly to date. *Blog* <https://www.pacb.com/blog/puerto-rican-genome/> (2018).
- Paten, B. et al. Cactus: Algorithms for genome multiple sequence alignment. *Genome Res.* **21**, 1512–1528 (2011).
- Kurtz, S. et al. Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
- Dot: an interactive dot plot viewer for comparative genomics (GitHub, 2020); <https://github.com/dnanexus/dot> (GitHub, 2020).
- Kroenenberg, Z. N. et al. High-resolution comparative analysis of great ape genomes. *Science* **360**, pii: eaar6343 (2018).
- Vollger, M. R. et al. Improved assembly and variant detection of a haploid human genome using single-molecule, high-fidelity long reads. *Ann. Hum. Genet.* **84**, 125–140 (2020).
- Rajaraman, A. & Ullman, J. D. *Mining of Massive Datasets* (Cambridge Univ. Press, 2011).
- Myers, E. W. The fragment assembly string graph. *Bioinformatics* **21** (Suppl. 2), ii79–ii85 (2005).
- Anderson, R. J. & Wont, H. Wait-free parallel algorithms for the union-find problem. In *Proc. Annual ACM Symposium on Theory of Computing Part F130073*, 370–380 (ACM, 1991).
- Onodera, T., Sadakane, K. & Shibuya, T. in *Algorithms in Bioinformatics* (eds Darling, A. & Stoye, J.) 338–348 (Springer, 2013).
- Chin, C. S. et al. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat. Methods* **10**, 563–569 (2013).
- Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat. Methods* **12**, 733–735 (2015).
- Paten, B., Herrero, J., Beal, K. & Birney, E. Sequence progressive alignment, a framework for practical large-scale probabilistic consistency alignment. *Bioinformatics* **25**, 295–301 (2008).
- Wright (ONT) provided advice for Medaka. D. Garalde and R. Dokos (ONT) provided advice on the PromethION for parallelized DNA sequencing and basecalling. We are grateful to AWS for hosting the data via their AWS Public Dataset Program. A.P. and S.K. were supported by the Intramural Research Program of the National Human Genome Research Institute, National Institutes of Health. This work used the computational resources of the NIH HPC Biowulf cluster (<https://hpc.nih.gov>). S. Bell and C. Weaver from Chan Zuckerberg Initiative (CZI) provided support on development and documentation. The CZI further supported this effort by funding the usage of AWS for the project. Certain commercial equipment, instruments, or materials are identified to specify adequately experimental conditions or reported results. Such identification does not imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the equipment, instruments or materials identified are necessarily the best available for the purpose. This work was supported, in part, by the National Institutes of Health (award nos. 2U41HG007234, 5U54HG007990, 5T32HG008345-04, R01HG010053 and U01HL137183 to B.P. and D.H.; R01HG010329 to S.R.S. and D.H.; 3U24HG009084-03S1 and 5R03HG009730-03 to H.E.O.), by Oxford Nanopore Research grant no. SC20130149 (M.A.), the Howard Hughes Medical Institute (D.H.), W.M. Keck Foundation (grant no. DT06172015) and the National Institute of Standards and Technology (grant no. 70NANB18H224).

Author contributions

B.P., M.J. and P.C. designed and executed the development of the study. P.C. is the core developer of Shasta. K.S., T.P., R.L.-R., M.H. and H.E.O. contributed equally in the core analysis presented in this study. T.P. and B.P. developed MarginPolish and K.S. developed HELEN. R.L. developed models for Shasta and MarginPolish and performed the MHC analysis. M.H. performed scaffolding using HiC data. M.J. and H.E.O. developed methods for nanopore long-read sequencing, performed experiments, basecalling and data wrangling. C.B. helped with MHC analysis. J.A. performed CAT analysis. K.T. performed cell culture. N.M. prepared HiC libraries. S.K. and A.P. performed trio-binning and Canu assemblies and helped with assembly evaluation. B.P., A.P. and F.J.S. helped with sample selection for sequencing. S.M. and V.C. helped with protocol development and nanopore sequencing. J.M.Z. helped with assembly evaluation and provided GIAB truth sets. K.J.L. and D.K. provided early access and help with size selection using Circulomics Short-Read Eliminator kits. E.E.E., M.R.V., T.M., M.S. and K.M.M. helped with assembly evaluation and provided BAC libraries. E.G. and J.M. performed structural variation analysis. S.S., D.H., R.E.G., M.A. and K.H.M. provided scientific feedback on the analyses. All authors contributed to writing and editing the manuscript.

Competing interests

M.A. is a paid consultant to ONT. V.C. and S.M. are employees of ONT.

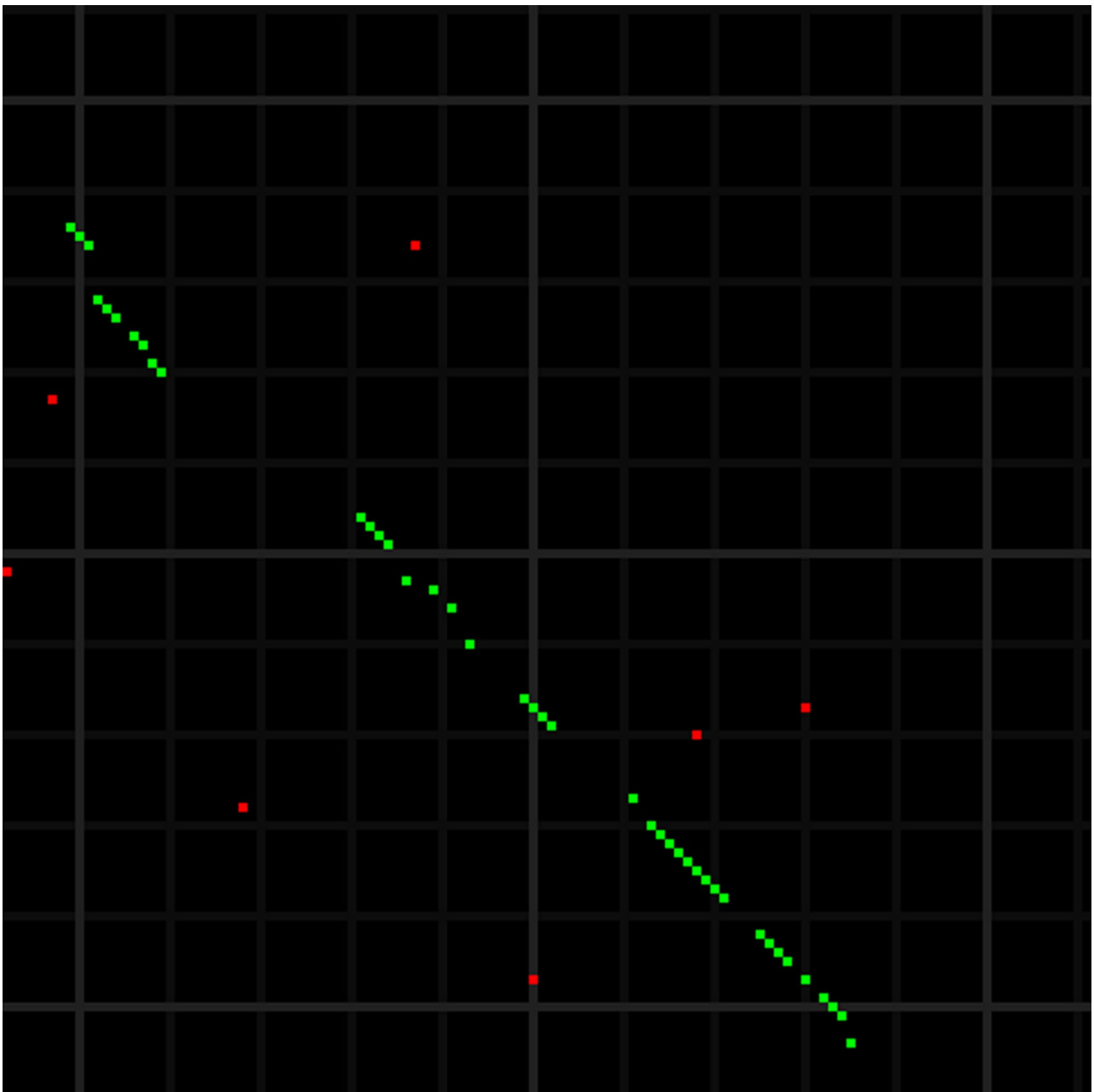
Additional information

Extended data is available for this paper at <https://doi.org/10.1038/s41587-020-0503-6>.

Supplementary information is available for this paper at <https://doi.org/10.1038/s41587-020-0503-6>.

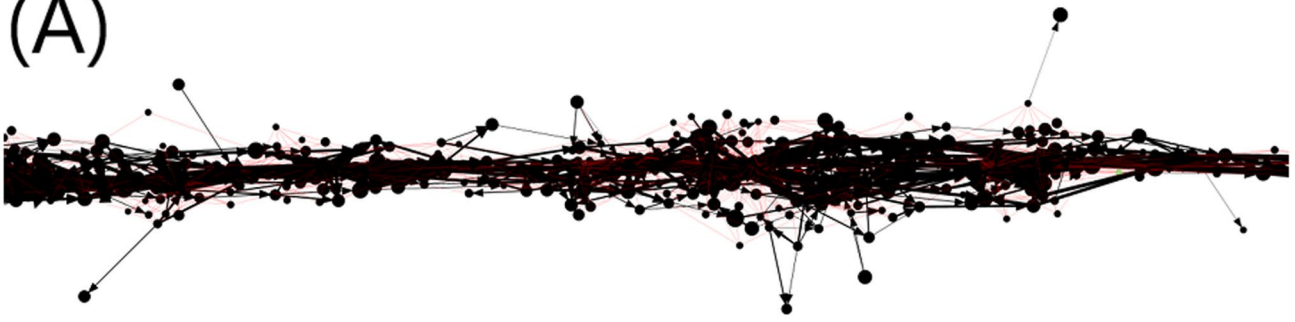
Correspondence and requests for materials should be addressed to P.C., M.J. or B.P.

Reprints and permissions information is available at www.nature.com/reprints.

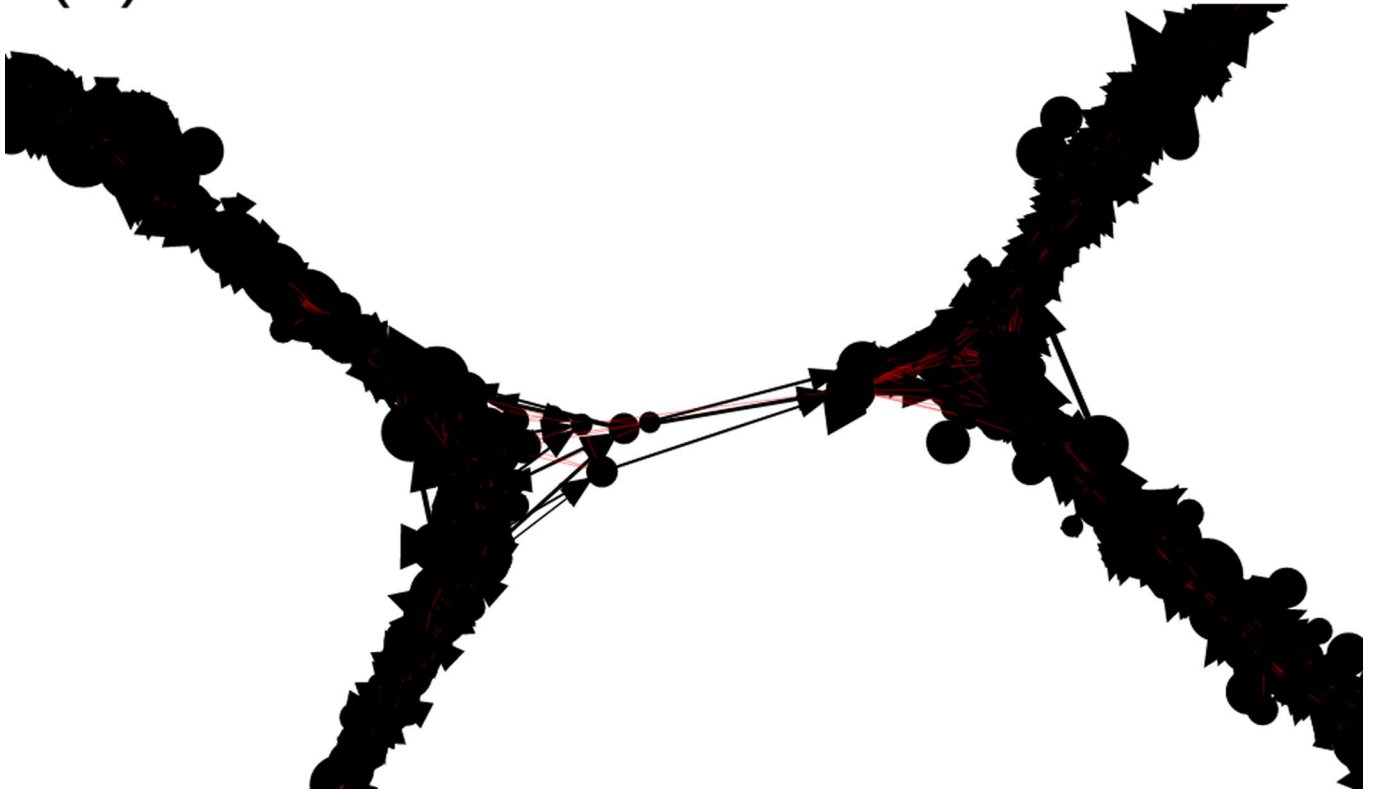


Extended Data Fig. 2 | Marker Alignment. A marker alignment represented as a dot-plot. Elements that are identical between the two sequences are displayed in green or red - the ones in green are the ones that are part of the optimal alignment computed by the Shasta assembler. Because of the much larger alphabet, matrix elements that are identical between the sequences but are not part of the optimal alignment are infrequent. Each alignment matrix element here corresponds on average to a 13 13 block in the alignment matrix in raw base sequence.

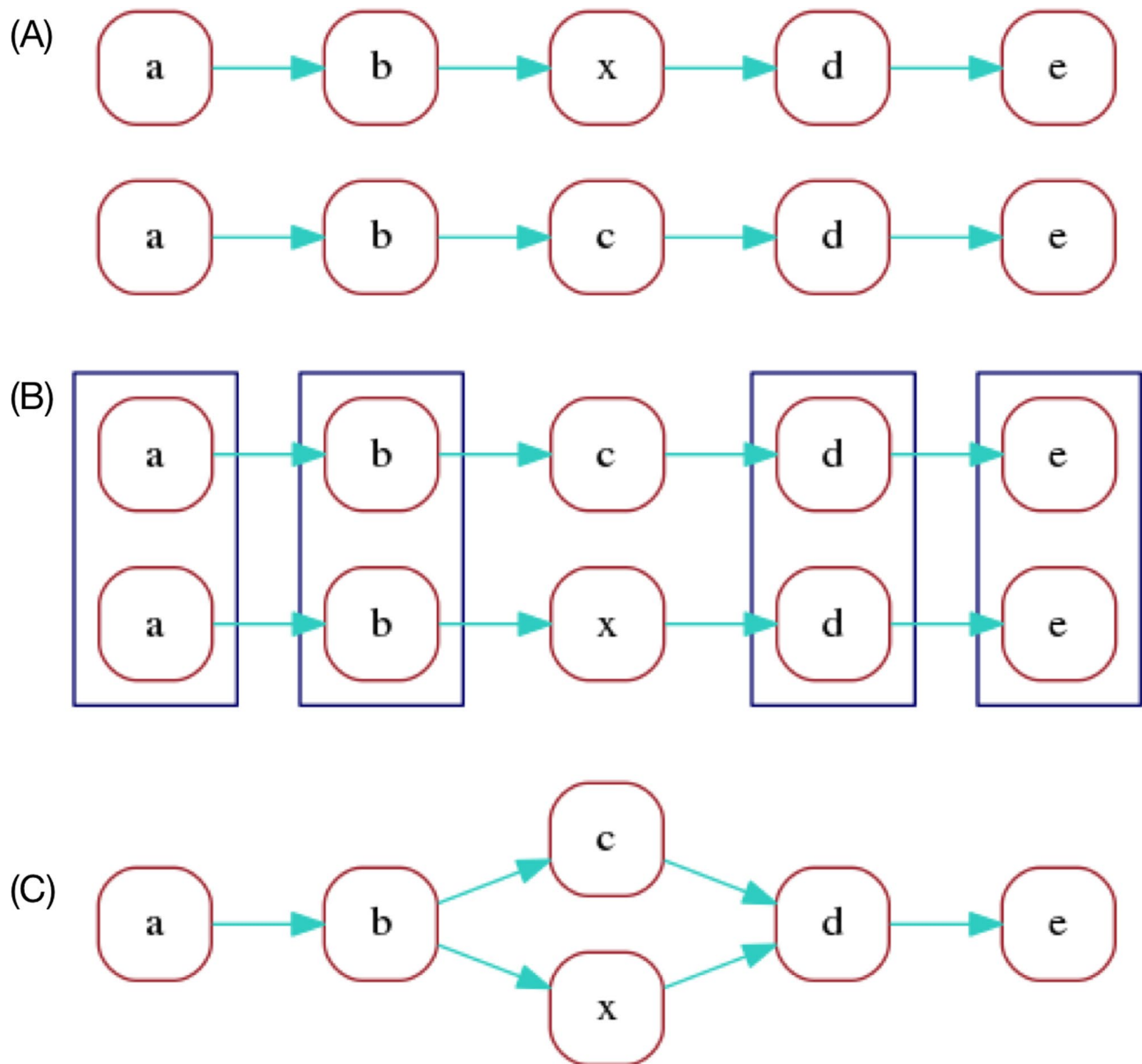
(A)



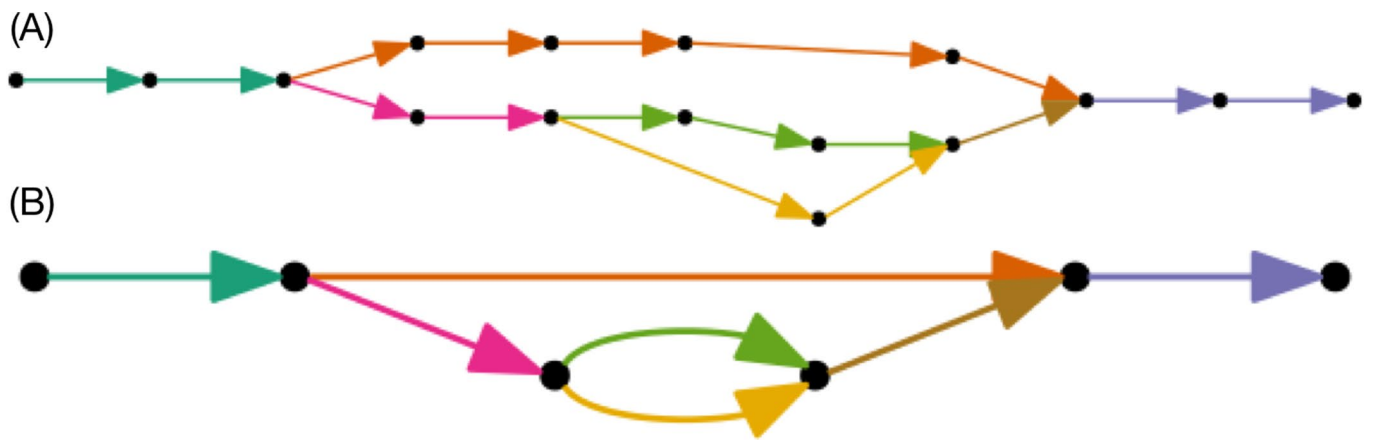
(B)



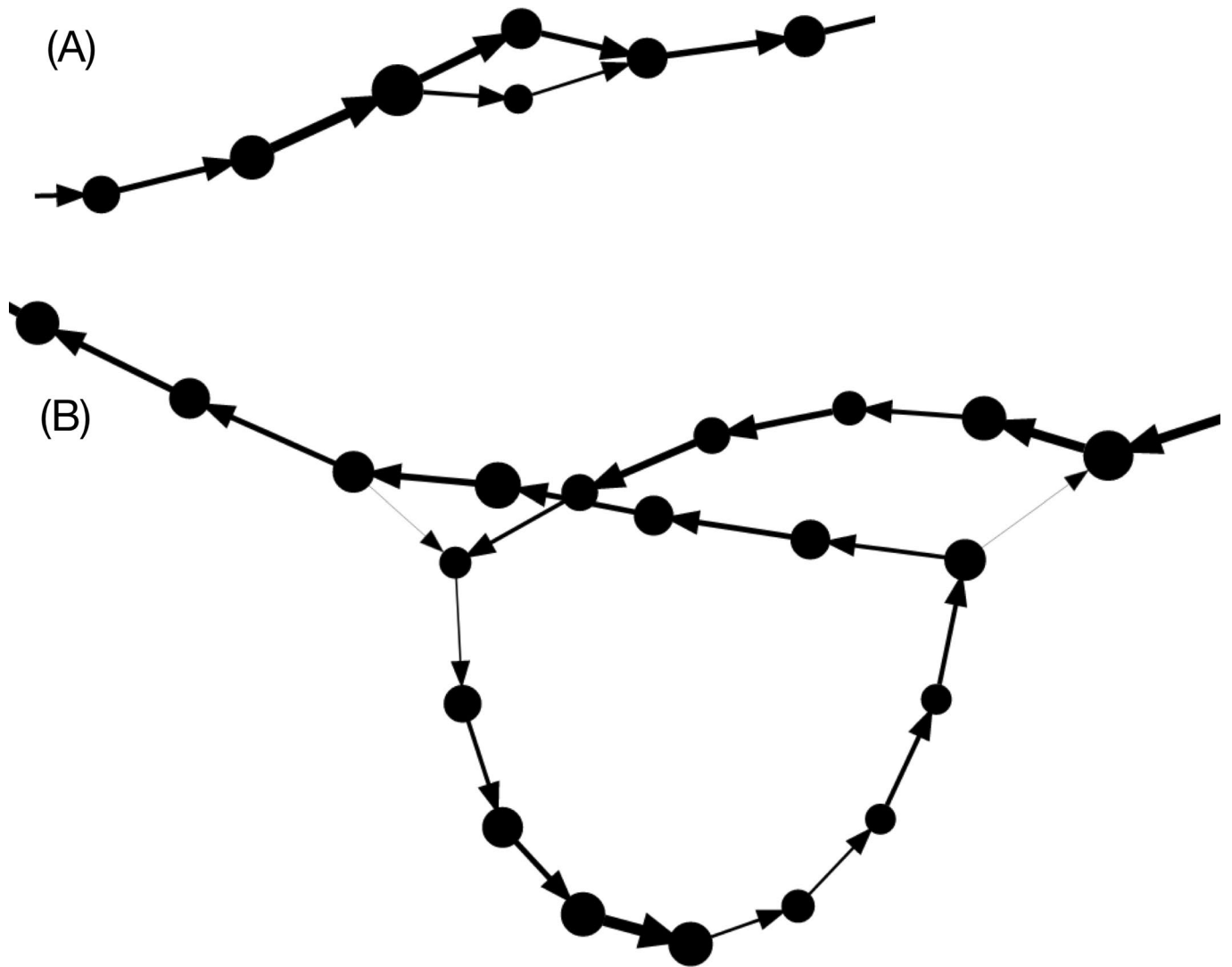
Extended Data Fig. 3 | Read Graph. An example of a portion of the read graph (A) as displayed by the Shasta http server, and (B) showing obviously incorrect connections.



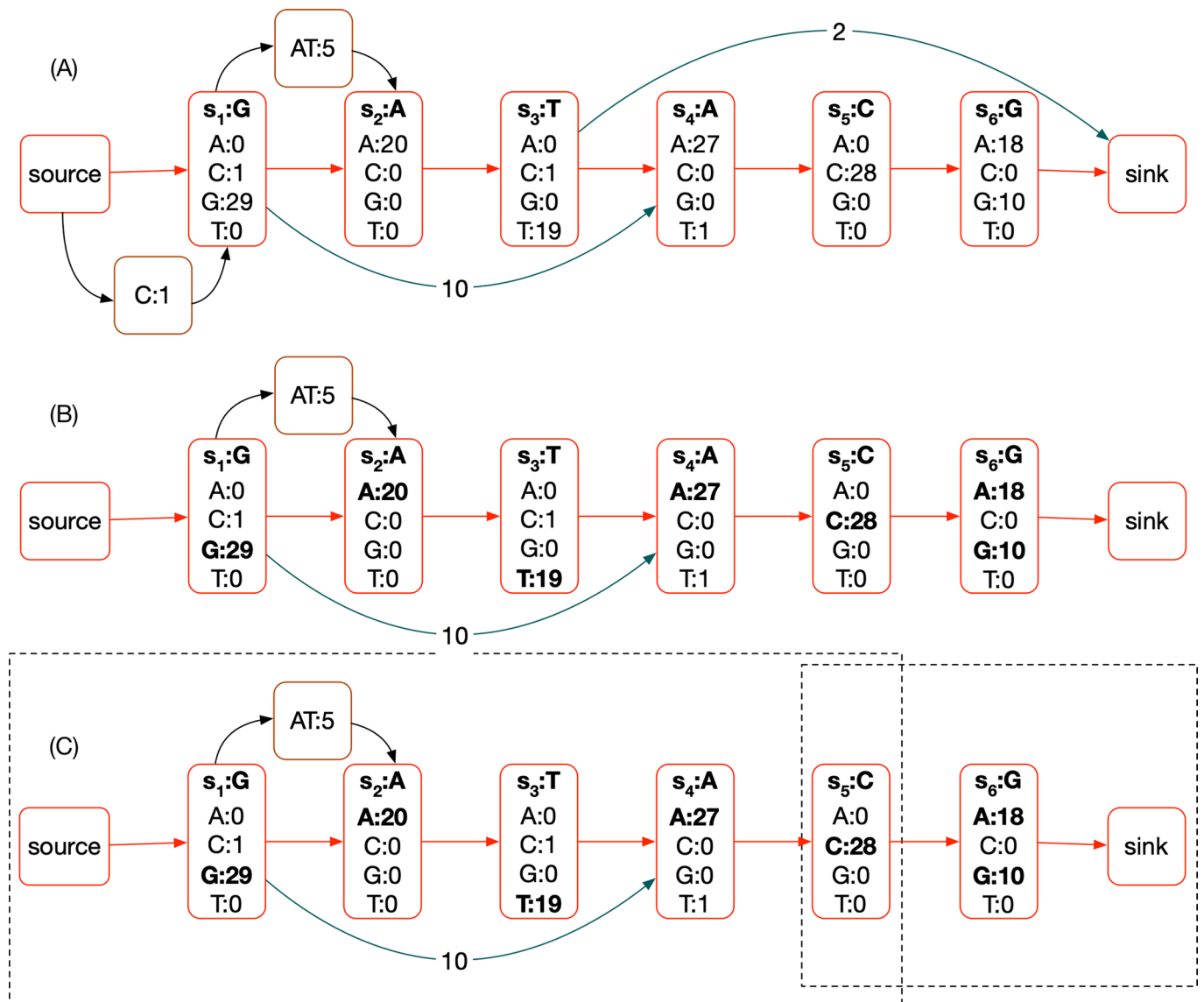
Extended Data Fig. 4 | Marker Graph. An illustration of marker graph construction for two sequences.



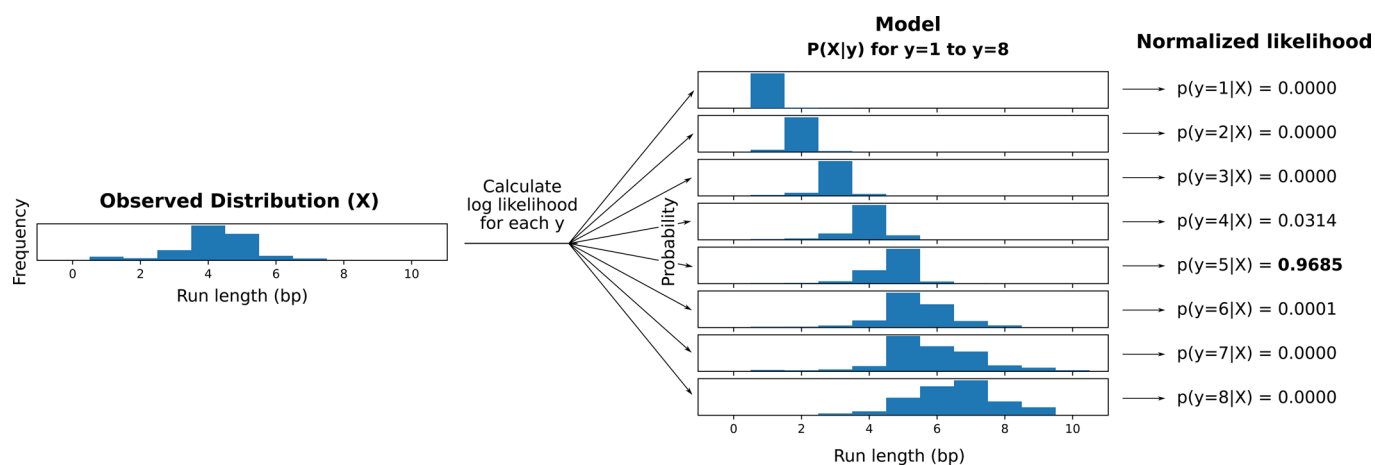
Extended Data Fig. 5 | Assembly Graph. (A) A marker graph with linear sequence of edges colored. (B) The corresponding assembly graph. Colors were chosen to indicate the correspondence to marker graph edges.



Extended Data Fig. 6 | Bubbles. (A) A simple bubble. (B) A superbubble.



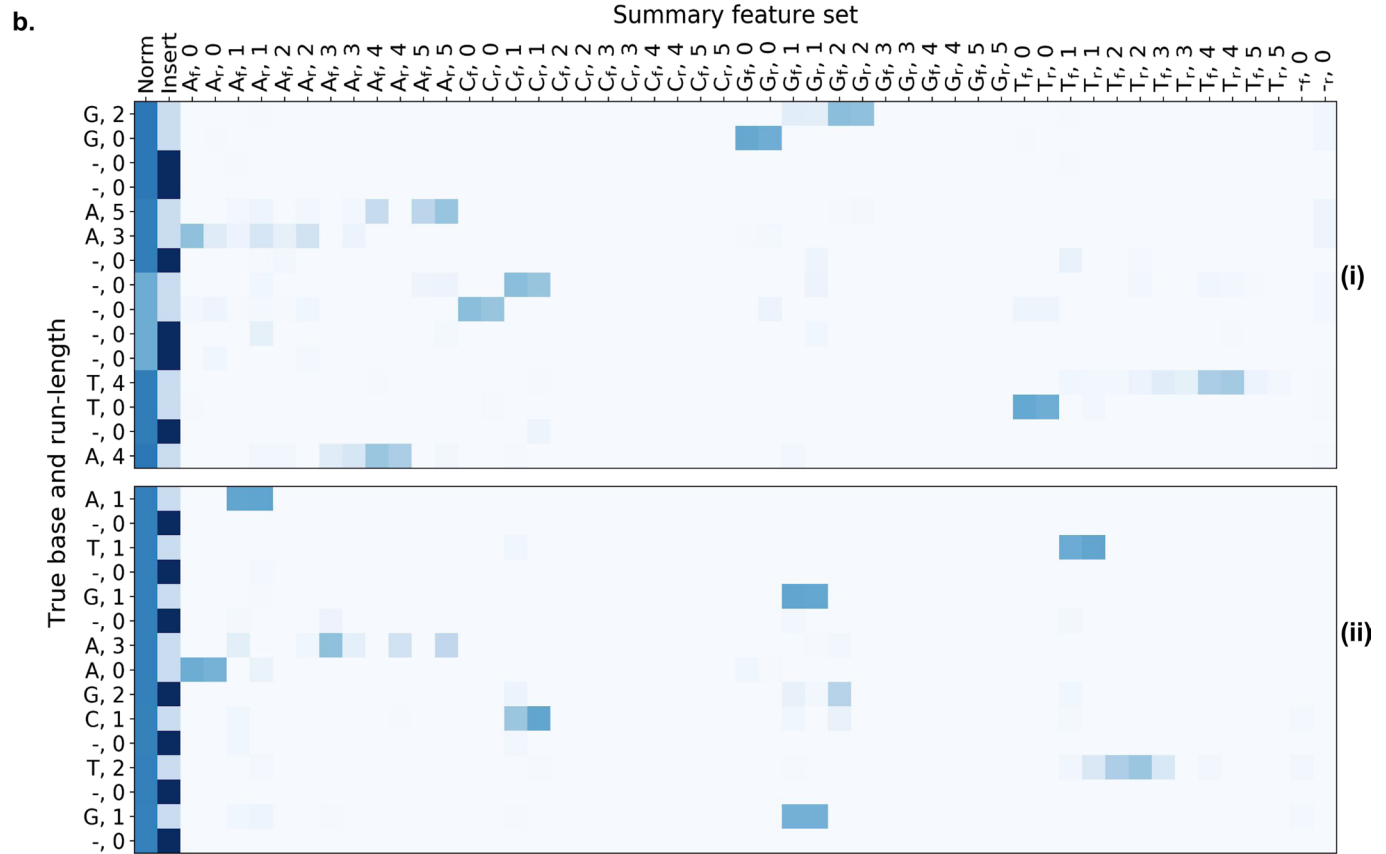
Extended Data Fig. 7 | POA Example. (A) An example POA, assuming approximately 30x read coverage. The backbone is shown in red. Each non-source/sink node has a vector of weights, one for each possible base. Deletion edges are shown in teal, they also each have a weight. Finally insertion nodes are shown in brown, each also has a weight. (B) A pruned POA, removing deletions and insertions that have less than a threshold weight and highlighting plausible bases in bold. There are six plausible nucleotide sequences represented by paths through the POA and selections of plausible base labels: G;AT;A;T;A;C:A, G;AT;A;T;A;C:G, G;A;T;A;C:A, G;A;T;A;C:G, G;A;C:A, G;A;C:G. To avoid the combinatorial explosion of such enumeration we identify subgraphs (C) and locally enumerate the possible subsequences in these regions independently (dotted rectangles identify subgraphs selected). In each subgraph there is a source and sink node that does not overlap any proposed edit.



Extended Data Fig. 8 | RLE Inference Distributions. Visual representation of run length inference. This diagram shows how a consensus run length is inferred for a set of aligned lengths (X) that pertain to a single position. The lengths are factored and then iterated over, and log likelihood is calculated for every possible true length up to a predefined limit. Note that in this example, the most frequent observation (4bp) is not the most likely true length (5bp) given the model.

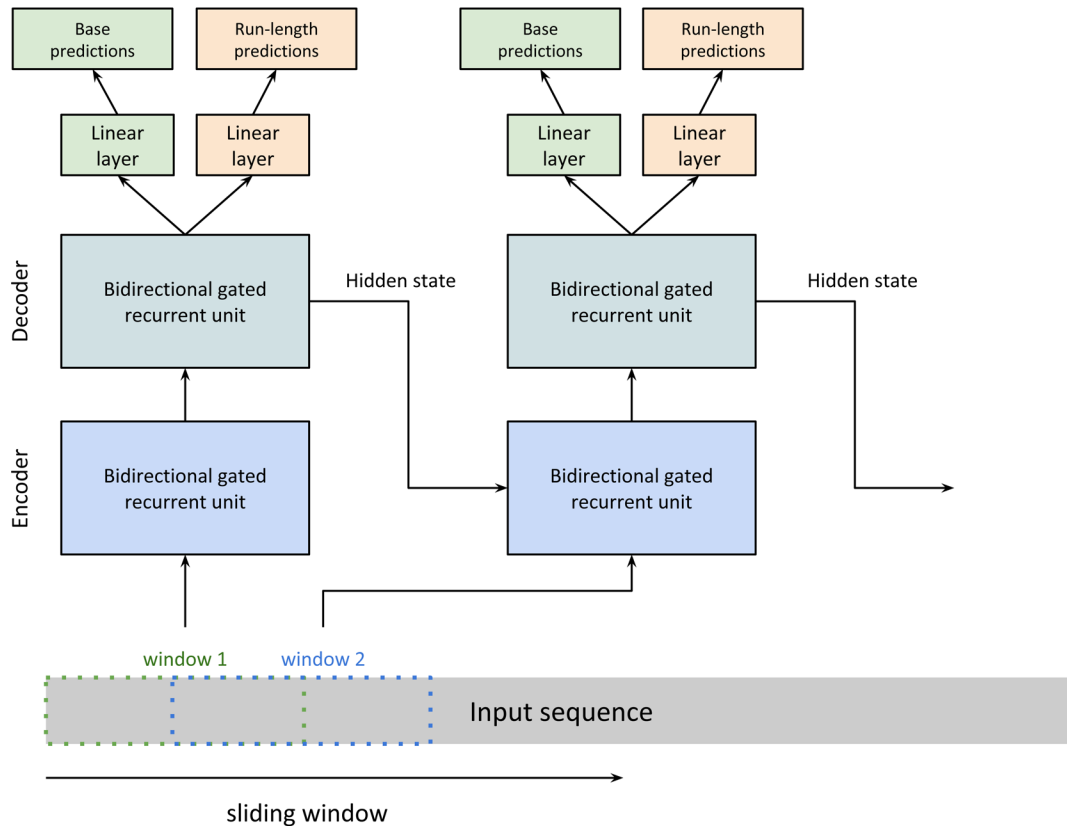
a.

<p>(i) Assembly sequence: GGAAAAAAAAACATTTTAAAA True sequence: GGAAAAAAAA - - TTTTAAAA</p> <p>Assembly sequence in run-length: G A A C A T A 2 5 3 1 1 4 4</p> <p>Truth sequence in run-length: G A A - - T A 2 5 3 0 0 4 4</p>	<p>(ii) Assembly sequence: ATGAAA - - CTTG True sequence: ATGAAAGGCTTG</p> <p>Assembly sequence in run-length: A T G A C T G 1 1 1 3 1 2 1</p> <p>Truth sequence in run-length: A T G A G C T G 1 1 1 3 2 1 2 1</p>
--	--



Extended Data Fig. 9 | MarginPolish HELEN Image Generation. A graphical representation of images from two labeled regions selected to demonstrate: the encoding of a single POA node into two run-length blocks (i), a true deletion (i), and a true insert (ii). (a) shows the alignment in raw and run-length space, (b) shows the features as they are exported to HELEN. The y-axis shows truth labels for nucleotides and run-lengths, the x-axis describes features in the images, and colors show associated weights.

Multi-task learning with hard parameter sharing



Extended Data Fig. 10 | HELEN Model. The sequence-to-sequence model implemented in Helen.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

We used Nanopore Technologies MinKnow software that comes with the sequencing device by default. The Guppy basecaller model we used was Guppy 2.3.5.

Data analysis

The online methods and supplementary notes sections describe all available software and the commands we used for data analysis.

Base calling: Oxford Nanopore Technology provided Guppy basecaller v2.3.5 (commercial product)
 Shasta: <https://github.com/chanzuckerberg/shasta> (MIT License)
 Canu: <https://github.com/marbl/canu> (GNU General Public License, version 2)
 WTDBG2: <https://github.com/ruanjue/wtdbg2> (GPL-3.0 License)
 Flye: <https://github.com/fenderglass/Flye> (BSD-3-Clause License)
 Racon 4x: https://github.com/rlorigro/nanopore_assembly_and_polishing_assessment (MIT License)
 Medaka: <https://github.com/nanoporetech/medaka> (MPL-2.0 License)
 Pomoxis (mini_align and assess_assembly): <https://github.com/nanoporetech/pomoxis> (MPL-2.0 License)
 Minimap2: <https://github.com/lh3/minimap2> (MIT License)
 MarginPolish: <https://github.com/UCSC-nanopore-cgl/MarginPolish> (MIT License)
 HELEN: <https://github.com/kishwarshafin/helen> (MIT License)
 CAT: <https://github.com/ComparativeGenomicsToolkit/Comparative-Annotation-Toolkit> (Apache-2.0 License)
 Read alignment identity: https://github.com/rlorigro/nanopore_assembly_and_polishing_assessment (MIT License)
 QUAST: <https://github.com/rlorigro/quast> (GNU General Public License, Version 2)
 misassembly stats: https://github.com/kishwarshafin/helen/blob/master/modules/python/helper/quast_sv_extractor.py (MIT License)
 Run-length confusion matrix: https://github.com/rlorigro/runlength_analysis/ (MIT License)
 Runtime and cost analysis: <https://github.com/rlorigro/TaskManager> (MIT License)
 BAC analysis: <https://github.com/skoren/bacValidation> (Public domain)

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

Sequence data including raw signal files (FAST5), base-calls (FASTQ), Illumina Hi-C data (FASTQ), are publicly hosted on Amazon Web Services Public Datasets program and available for download via GitHub here: <https://github.com/human-pangenomics/hpgp-data>

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

- Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No sample size calculation was performed. Cells were grown in batches and aliquoted into 50 million sized pellets. We chose 50M sized cell pellets to ensure we could isolate sufficient DNA for the experimental protocols. The 50M size was selected based on prior literature (Jain et al. NBT 2018).
Data exclusions	No data were excluded from the analyses.
Replication	We performed three replicate experiments for data generation. The replicates were consistent in data quality, as demonstrated in Figure 1.
Randomization	No randomization was performed by group; however, this was not relevant to our study as data from groups was pooled and analyzed together.
Blinding	Investigators were not blind to group allocation during experiments, as all data from groups was pooled together.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input type="checkbox"/>	<input checked="" type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Eukaryotic cell lines

Policy information about [cell lines](#)

Cell line source(s)	Coriell Institute. Cell lines: HG002, HG003, HG004, HG02055, HG02080, HG03492, HG00733, HG03098, HG01243, HG02723, HG01109
Authentication	The cell lines used were not authenticated.
Mycoplasma contamination	The cell lines used were not tested for mycoplasma contamination.
Commonly misidentified lines (See ICLAC register)	The cell lines used are not in the register of commonly misidentified lines.