









OPEN

Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain^{1,13} , Sergey Koren^{2,13}, Karen H Miga^{1,13}, Josh Quick^{3,13}, Arthur C Rand^{1,13}, Thomas A Sasani^{4,5,13} , John R Tyson^{6,13}, Andrew D Beggs⁷ , Alexander T Dilthey² , Ian T Fiddes¹, Sunir Malla⁸, Hannah Marriott⁸, Tom Nieto⁷, Justin O'Grady⁹ , Hugh E Olsen¹, Brent S Pedersen^{4,5}, Arang Rhie² , Hollian Richardson⁹, Aaron R Quinlan^{4,5,10} , Terrance P Snutch⁶, Louise Tee⁷, Benedict Paten¹, Adam M Phillippy², Jared T Simpson^{11,12}, Nicholas J Loman³ & Matthew Loose⁸ 

We report the sequencing and assembly of a reference genome for the human GM12878 Utah/Ceph cell line using the MinION (Oxford Nanopore Technologies) nanopore sequencer. 91.2 Gb of sequence data, representing ~30× theoretical coverage, were produced. Reference-based alignment enabled detection of large structural variants and epigenetic modifications. *De novo* assembly of nanopore reads alone yielded a contiguous assembly (NG50 ~3 Mb). We developed a protocol to generate ultra-long reads (N50 > 100 kb, read lengths up to 882 kb). Incorporating an additional 5× coverage of these ultra-long reads more than doubled the assembly contiguity (NG50 ~6.4 Mb). The final assembled genome was 2,867 million bases in size, covering 85.8% of the reference. Assembly accuracy, after incorporating complementary short-read sequencing data, exceeded 99.8%. Ultra-long reads enabled assembly and phasing of the 4-Mb major histocompatibility complex (MHC) locus in its entirety, measurement of telomere repeat length, and closure of gaps in the reference human genome assembly GRCh38.

The human genome is used as a yardstick to assess performance of DNA sequencing instruments^{1–5}. Despite improvements in sequencing technology, assembling human genomes with high accuracy and completeness remains challenging. This is due to size (~3.1 Gb), heterozygosity, regions of GC% bias, diverse repeat families, and segmental duplications (up to 1.7 Mbp in size) that make up at least 50% of the genome⁶. Even more challenging are the pericentromeric, centromeric, and acrocentric short arms of chromosomes, which contain satellite DNA and tandem repeats of 3–10 Mb in length^{7,8}. Repetitive structures pose challenges for *de novo* assembly using “short read” sequencing technologies, such as Illumina's. Such data, while enabling highly accurate genotyping in non-repetitive regions, do not provide contiguous *de novo* assemblies. This limits the ability to reconstruct repetitive sequences, detect complex structural variation, and fully characterize the human genome.

Single-molecule sequencers, such as Pacific Biosciences' (PacBio), can produce read lengths of 10 kb or more, which makes *de novo* human genome assembly more tractable⁹. However, single-molecule sequencing reads have significantly higher error rates compared with Illumina sequencing. This has necessitated development of *de novo* assembly

algorithms and the use of long noisy data in conjunction with accurate short reads to produce high-quality reference genomes¹⁰. In May 2014, the MinION nanopore sequencer was made available to early-access users¹¹. Initially, the MinION nanopore sequencer was used to sequence and assemble microbial genomes or PCR products^{12–14} because the output was limited to 500 Mb to 2 Gb of sequenced bases. More recently, assemblies of eukaryotic genomes including yeasts, fungi, and *Caenorhabditis elegans* have been reported^{15–17}.

Recent improvements to the protein pore (a laboratory-evolved *Escherichia coli* CsgG mutant named R9.4), library preparation techniques (1D ligation and 1D rapid), sequencing speed (450 bases/s), and control software have increased throughput, so we hypothesized that whole-genome sequencing (WGS) of a human genome might be feasible using only a MinION nanopore sequencer^{17–19}.

We report sequencing and assembly of a reference human genome for GM12878 from the Utah/CEPH pedigree, using MinION R9.4 1D chemistry, including ultra-long reads up to 882 kb in length. GM12878 has been sequenced on a wide variety of platforms, and has well-validated variation call sets, which enabled us to benchmark our results²⁰.

¹UC Santa Cruz Genomics Institute, University of California, Santa Cruz, California, USA. ²Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, Bethesda, Maryland, USA. ³Institute of Microbiology and Infection, University of Birmingham, Birmingham, UK. ⁴Department of Human Genetics, University of Utah, Salt Lake City, Utah, USA. ⁵USTAR Center for Genetic Discovery, University of Utah, Salt Lake City, Utah, USA. ⁶Michael Smith Laboratories and Djeval Mowafaghian Centre for Brain Health, University of British Columbia, Vancouver, Canada. ⁷Surgical Research Laboratory, Institute of Cancer & Genomic Science, University of Birmingham, UK. ⁸DeepSeq, School of Life Sciences, University of Nottingham, UK. ⁹Norwich Medical School, University of East Anglia, Norwich, UK. ¹⁰Department of Biomedical Informatics, University of Utah, Salt Lake City, Utah, USA. ¹¹Ontario Institute for Cancer Research, Toronto, Canada. ¹²Department of Computer Science, University of Toronto, Toronto, Canada. ¹³These authors contributed equally to this work. Correspondence should be addressed to N.J.L. (n.j.loman@bham.ac.uk) or M.L. (matt.loose@nottingham.ac.uk).

Received 20 April 2017; accepted 11 December 2017; published online 29 January 2018; doi:10.1038/nbt.4060

RESULTS

Sequencing data set

Five laboratories collaborated to sequence DNA from the GM12878 human cell line. DNA was sequenced directly (avoiding PCR), thus preserving epigenetic modifications such as DNA methylation. 39 MinION flow cells generated 14,183,584 base-called reads containing 91,240,120,433 bases with a read N50 (the read length such that reads of this length or greater sum to at least half the total bases) of 10,589 bp (Supplementary Tables 1–4). Ultra-long reads were produced using 14 additional flow cells. Read lengths were longer when the input DNA was freshly extracted from cells compared with using Coriell-supplied DNA (Fig. 1a). Average yield per flow cell (2.3 Gb) was unrelated to DNA preparation methods (Fig. 1b). 94.15% of reads had at least one alignment to the human reference (GRCh38) and 74.49% had a single alignment over 90% of their length. Median coverage depth was 26-fold, and 96.95% (3.01/3.10 Gbp) of bases of the reference were covered by at least one read (Fig. 1c). The median identity of reads was 84.06% (82.73% mean, 5.37% s.d.). No length bias was observed in the error rate with the MinION (Fig. 1d).

Base-caller evaluation

The base-calling algorithm used to decode raw ionic current signal can affect sequence calls. To analyze this effect we used reads mapping

to chromosome 20 and compared base-calling with Metrichor (an LSTM-RNN base-caller) and Scrapie, an open-source transducer neural network (Online Methods). Of note, we observed that a fraction of the Scrapie output (4.7% reads, 14% bases) was composed of low-complexity sequence (Supplementary Fig. 1), which we removed before downstream analysis.

To assess read accuracy we realigned reads from each base-caller using a trained alignment model²¹. Alignments generated by the Burrows–Wheeler Aligner Maximal Exact Matches (BWA-MEM) were chained such that each read had at most one maximal alignment to the reference sequence (scored by length). The chained alignments were used to derive the maximum likelihood estimate of alignment model parameters²², and the trained model used to realign the reads. The median identity after realignment for Metrichor was 82.43% and for Scrapie, 86.05%. We observed a purine-to-purine substitution bias in chained alignments where the model was not used (Supplementary Fig. 2). The alignments produced by the trained model showed an improved substitution error rate, decreasing the overall transversion rate, but transition errors remained dominant.

To measure potential bias at the *k*-mer level, we compared counts of 5-mers in reads derived from chromosome 20. In Metrichor reads, the most underrepresented 5-mers were A/T-rich homopolymers. The most overrepresented *k*-mers were G/C-rich and non-homopolymeric

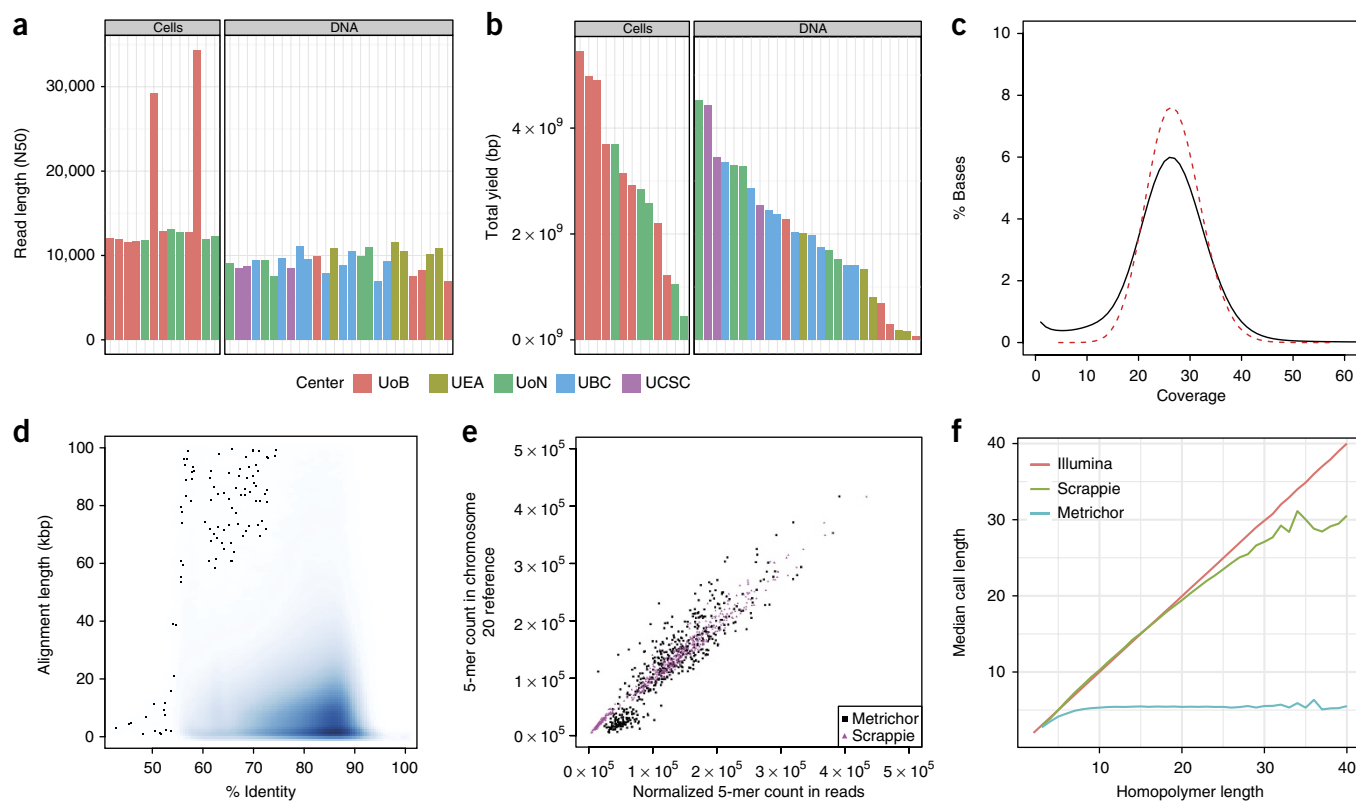


Figure 1 Summary of data set. (a) Read length N50s by flow cell, colored by sequencing center. Cells: DNA extracted directly from cell culture. DNA: pre-extracted DNA purchased from Coriell. UoB, Univ. Birmingham; UEA, Univ. East Anglia; UoN, Univ. Nottingham; UBC, Univ. British Columbia; UCSC, Univ. California, Santa Cruz. (b) Total yield per flow cell grouped as in a. (c) Coverage (black line) of GRCh38 reference compared to a Poisson distribution. The depth of coverage of each reference position was tabulated using samtools depth and compared with a Poisson distribution with $\lambda = 27.4$ (dashed red line). (d) Alignment identity compared to alignment length. No length bias was observed, with long alignments having the same identity as short ones. (e) Correlation between 5-mer counts in reads compared to expected counts in the chromosome 20 reference. (f) Chromosome 20 homopolymer length versus median homopolymer base-call length measured from individual Illumina and nanopore reads (Scrapie and Metrichor). Metrichor fails to produce homopolymer runs longer than ~5 bp. Scrapie shows better correlation for longer homopolymer runs, but tends to overcall short homopolymers (between 5 and 15 bp) and undercall long homopolymers (>15 bp). Plot noise for longer homopolymers is due to fewer samples available at that length.

Table 1 Summary of assembly statistics

| Assembly | Polishing | Contigs | No. bases (Mbp) | Max contig (kb) | NG50 (kb) | GRCh38 identity (%) | GM12878 identity (%) |
|------------------|-----------------|---------|-----------------|-----------------|-----------|---------------------|----------------------|
| WGS Metrichor | N/A | 2,886 | 2,646.01 | 27,160 | 2,964 | 95.20 | 95.74 |
| | Pilon x2 | | 2,763.18 | 28,413 | 3,206 | 99.29 | 99.88 |
| Chr 20 Metrichor | N/A | 85 | 57.83 | 7,393 | 3,047 | 94.90 | 95.50 |
| | Nanopolish | | 60.35 | 7,667 | 5,394 | 98.84 | 99.24 |
| | Pilon x2 | | 60.58 | 7,680 | 5,423 | 99.33 | 99.89 |
| | Nano + Pilon x2 | | 60.76 | 7,698 | 5,435 | 99.64 | 99.95 |
| Chr 20 Scrapie | N/A | 74 | 59.39 | 8,415 | 2,643 | 97.43 | 97.80 |
| | Nanopolish | | 60.15 | 8,521 | 2,681 | 99.12 | 99.44 |
| | Pilon x2 | | 60.36 | 8,541 | 2,691 | 99.64 | 99.95 |
| | Nano + Pilon x2 | | 60.34 | 8,545 | 2,691 | 99.70 | 99.96 |

Summary of assembly statistics. Whole genome assembly (WGA) was performed with reads base-called by Metrichor. Chromosome 20 was assembled with reads produced by Metrichor and Scrapie. All data sets contained 30× coverage of the genome/chromosome. The GRCh38 identities were computed based on 1-1 alignments to the GRCh38 reference including alt sites. A GM12878 reference was estimated using an independent sequencing data set²⁰.

(Supplementary Table 5). By contrast, Scrapie showed no underrepresentation of homopolymeric 5-mers and had a slight overrepresentation of A/T homopolymers. Overall, Scrapie showed the lowest *k*-mer representation bias (Fig. 1e). The improved homopolymer resolution of Scrapie was confirmed by inspection of chromosome 20 homopolymer calls versus the human reference (Fig. 1f and Supplementary Fig. 3)²³. Despite this reduced bias, whole-genome assembly and analyses proceeded with Metrichor reads, since Scrapie was still in early development at the time of writing.

De novo assembly of nanopore reads

We carried out a *de novo* assembly of the 30× data set with Canu²⁴ (Table 1). This assembly comprised 2,886 contigs with an NG50 contig size of 3 Mbp (NG50, the longest contig such that contigs of this length or greater sum to at least half the haploid genome size). The identity to GRCh38 was estimated as 95.20%. Canu was four-fold slower on the Nanopore data compared to a random subset of equivalent coverage of PacBio data requiring ~62K CPU hours. The time taken by Canu increased when the input was nanopore sequence reads because of systematic error in the raw sequencing data leading to reduced accuracy of the Canu-corrected reads, an intermediate output of the assembler. Corrected PacBio reads are typically >99% identical to the reference; our reads averaged 92% identity to the reference after correction (Supplementary Fig. 1b).

We aligned assembled contigs to the GRCh38 reference and found that our assembly was in agreement with previous GM12878 assemblies (Supplementary Fig. 4)²⁵. The number of structural differences (899) that we identified between GM12878 and GRCh38 was similar to that of a previously published PacBio assembly of GM12878 (692) and of other human genome assemblies^{5,24}, but with a higher than expected number of deletions, due to consistent truncation of homopolymer and low-complexity regions (Supplementary Fig. 5 and Supplementary Table 6). Consensus identity of our assembly with GRCh38 was estimated to be 95.20% (Table 1). However, GRCh38 is a composite of multiple human haplotypes, so this is a lower bound on accuracy. Comparisons with independent Illumina data from GM12878 yielded a higher accuracy estimate of 95.74%.

Despite the low consensus accuracy, contiguity was good. For example, the assembly included a single ~3-Mbp contig that had all class I human leukocyte antigens (HLA) genes from the major histocompatibility complex (MHC) region on chromosome 6, a region notoriously difficult to assemble using short reads. The more repetitive class II HLA gene locus was fragmented but most genes were present in a single contig.

Genome polishing

To improve the accuracy of our assembly we mapped previously generated whole-genome Illumina data (SRA: ERP001229) to each contig using BWA-MEM and corrected errors using Pilon. This improved the estimated accuracy of our assembly to 99.29% versus GRCh8 and 99.88% versus independent GM12878 sequencing (Table 1 and Supplementary Fig. 6)²⁶. This estimate is a lower bound as true heterozygous variants and erroneously mapped sequences decrease identity. Recent PacBio assemblies of mammalian genomes that were assembled *de novo* and polished with Illumina data exceed 99.95%^{9,27}. Pilon cannot polish regions that have ambiguous short-read mappings, that is, in repeats. We also compared the accuracy of our polished assembly in regions with expected coverage versus those that had low-quality mappings (either lower coverage or higher than expected coverage with low mapping quality) versus GRCh38. When compared to GRCh38, accuracy in well-covered regions increased to 99.32% from the overall accuracy of 99.29%, while the poorly covered regions accuracy dropped to 98.65%.

For further evaluation of our assembly, we carried out comparative annotation before and after polishing (Supplementary Table 7). 58,338 genes (19,436 coding, 96.4% of genes in GENCODE V24, 98.2% of coding genes) were identified representing 179,038 transcripts in the polished assembly. Reflecting the assembly's high contiguity, only 857 (0.1%) of genes were found on two or more contigs.

Alternative approaches to improve assembly accuracy using different base-callers and exploiting the ionic current signal were attempted on a subset of reads from chromosome 20. Assembly consensus improvement using raw output is commonly used when assembling single-molecule data. To quantify the effect of base-calling on the assembly, we reassembled the read sets from Metrichor and Scrapie with the same Canu parameters used for the whole-genome data set. While all assemblies had similar contiguity, using Scrapie reads improved accuracy from 95.74% to 97.80%. Signal-level polishing of Scrapie-assembled reads using nanopolish increased accuracy to 99.44%, and polishing with Illumina data brought the accuracy up to 99.96% (Table 1).

Analysis of sequences not in the assembly

To investigate sequences omitted from the primary genome analysis, we assessed 1,425 contigs filtered from Canu due to low coverage, or contigs that were single reads with many shorter reads within (26 Mbp), or corrected reads not incorporated into contigs (10.4 Gbp). Most sequences represented repeat classes, for example, long interspersed nuclear elements (LINEs) and short interspersed nuclear elements (SINEs) (Supplementary Fig. 7), observed in similar

proportion in the primary assembly, with the exception of satellite DNAs known to be enriched in human centromeric regions. These satellites were enriched 2.93× in the unassembled data and 7.9× in the Canu-filtered contigs. We identified 56 assembled contigs containing centromere repeat sequences specific to each of the 22 autosomes and X chromosome. The largest assembled satellite in these contigs is a 94-kbp tandem repeat specific to centromere 15 (D15Z1, tig00007244).

SNP and SV genotyping

Using SVTyper²⁸ and Platinum Illumina WGS alignments, we genotyped 2,414 GM12878 structural variants (SVs), which were previously identified using LUMPY and validated with PacBio and/or Moleculo reads²⁹. We then genotyped the same SVs using alignments of our nanopore reads from the 30×-coverage data set and a modified version of SVTyper. We measured the concordance of genotypes at each site in the Illumina- and nanopore-derived data, deducing the sensitivity of SV genotyping as a function of nanopore sequencing depth (Fig. 2a). When all 39 flow cells were used, nanopore data recovered 91% of high-confidence SVs with a false-positive rate of 6%. Illumina and nanopore genotypes agreed at 81% of heterozygous sites and 90% of homozygous alternate sites. Genotyping heterozygous SVs using nanopore alignments was limited when homopolymer stretches occur at the breakpoints of these variants (Supplementary Fig. 8a). We determined Illumina, nanopore, and PacBio genotype concordance at a set of 2,192 deletions common to our high-confidence set and a genotyped SV call set derived from PacBio sequencing of GM12878 (refs. 5,30). PacBio and Illumina genotypes agreed at 94% of heterozygous and 79% of homozygous alternate deletions; nanopore and Illumina genotypes agreed at 90% of heterozygous and 90% of homozygous alternate sites; nanopore and PacBio genotypes agreed at 91% of heterozygous and 76% of homozygous alternate sites. Nearly a quarter (44) of the homozygous alternate sites at which PacBio and Illumina genotypes disagreed overlapped SINEs or LINEs. By manual inspection in the integrative genomics viewer (IGV)³¹, we observed that sequencing reads were spuriously aligned at these loci and likely drove the discrepancy in predicted genotypes (Supplementary Fig. 8b).

We evaluated nanopore data for calling genotypes at known single-nucleotide polymorphisms (SNPs) using the ionic current by calling genotypes at non-singleton SNPs on chromosome 20 from phase 3 of the 1000 Genomes³² and comparing these calls to Illumina Platinum Genome calls (Fig. 2b). 99.16% of genotype calls were correct (778,412 out of 784,998 sites). This result is dominated by the large number of homozygous reference sites. If we assess accuracy by the fraction of correctly called variant sites (heterozygous or homozygous non-reference), the accuracy of our caller is 91.40% (50,814 out of 55,595), with the predominant error being miscalling sites labeled homozygous in the reference as heterozygous (3,217 errors). Genotype accuracy, when only considering sites annotated as variants in the platinum call set, is 94.83% (50,814 correct out of 53,582).

Detection of epigenetic 5-methylcytosine modification

Changes in the ionic current when modified and unmodified bases pass through the MinION nanopores enable detection of epigenetic marks^{33,34}. We used nanopolish and SignalAlign to map 5-methylcytosine at CpG dinucleotides as detected in our sequencing reads against chromosome 20 of the GRCh38 reference^{35,36}. Nanopolish outputs a frequency of reads calling a methylated cytosine, and SignalAlign outputs a marginal probability of methylation summed over reads. We compared the output of both methods to published bisulfite

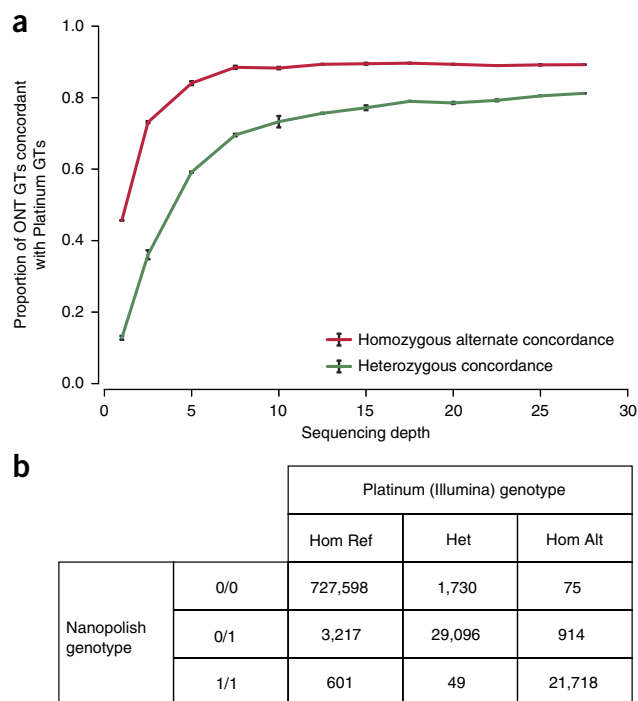


Figure 2 Structural variation and SNP genotyping. **(a)** Structural variant genotyping sensitivity using Oxford Nanopore Technologies (ONT) reads. Genotypes (GTs) were inferred for a set of 2,414 SVs using both Oxford Nanopore and Platinum Genomes (Illumina) alignments. Using alignments randomly subsampled to a given sequencing depth ($n = 3$), sensitivity was calculated as the proportion of ONT-derived genotypes that were concordant with Illumina-derived genotypes. **(b)** Confusion matrix for genotype-calling evaluation. Each cell contains the number of 1000 Genome sites for a particular nanopolish/platinum genotype combination.

sequencing data from the same DNA region (ENCFF835NTE). Good concordance of our data with the published bisulfite sequencing was observed; the r -values for nanopolish and SignalAlign were 0.895 and 0.779, respectively (Fig. 3 and Supplementary Figs. 9 and 10).

Ultra-long reads improve phasing and assembly contiguity

We modeled the contribution of read length to assembly quality, predicting that ultra-long read data sets ($N_{50} > 100$ kb) would substantially improve assembly contiguity (Fig. 4a). We developed a method to produce ultra-long reads by saturating the Oxford Nanopore Rapid Kit with high molecular weight DNA. In so doing we generated an additional 5× coverage (Supplementary Fig. 11). Two additional standard protocol flow cells generated a further 2× coverage and were used as controls for software and base-caller versions. The N_{50} read length of the ultra-long data set was 99.7 kb (Fig. 4b). Reads were impossible to align efficiently at first, because aligner algorithms are optimized for short reads. Further, CIGAR strings generated by ultra-long reads do not fit in the BAM format specification, necessitating the use of SAM or CRAM formats only (<https://github.com/samtools/hts-specs/issues/40>). Instead, we used GraphMap³⁷ to align ultra-long reads to GRCh38, which took >25K CPU hours (Supplementary Table 8). Software optimized for long reads, including NGM-LR³⁸ and Minimap2 (ref. 39), were faster: Minimap2 took 60 CPU hours. More than 80% of bases were in sequences aligned over 90% of their length with GraphMap and more than 60% with minimap2. Median alignment identity was 81% (83 with minimap2), slightly lower than observed for the control flow cells (83.46%/84.64%) and the original

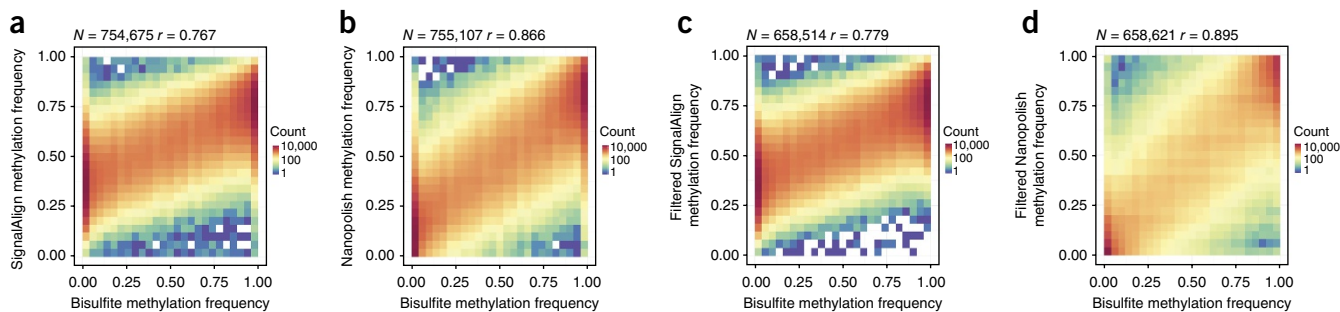


Figure 3 Methylation detection using signal-based methods. **(a)** SignalAlign methylation probabilities compared to bisulfite sequencing frequencies at all called sites. **(b)** Nanopolish methylation frequencies compared to bisulfite sequencing at all called sites. **(c)** SignalAlign methylation probabilities compared to bisulfite sequencing frequencies at sites covered by at least ten reads in the nanopore and bisulfite data sets; reads were not filtered for quality. **(d)** Nanopolish methylation frequencies compared to bisulfite sequencing at sites covered by at least ten reads in the nanopore and bisulfite data sets. A minimum log-likelihood threshold of 2.5 was applied to remove ambiguous reads. N = sample size, r = Pearson correlation coefficient.

data set (83.11%/84.32%). The longest full-length mapped read in the data set (aligned with GraphMap) was 882 kb, corresponding to a reference span of 993 kb.

The addition of 5 \times coverage ultra-long reads more than doubled the previous assembly NG50 to 6.4 Mbp and resolved the MHC locus into a single contig (Fig. 4c). In comparison, a 50 \times PacBio GM12878 data set with average read length of 4.5 kb assembled with an NG50 contig size of 0.9 Mbp⁵. Newer PacBio assemblies of a human haploid cell line, with mean read lengths greater than 10 kb, have reached contig NG50s exceeding 20 Mbp at 60 \times coverage²⁵. We subsampled this data set to a depth equivalent to ours (35 \times) and assembled, resulting in an NG50 of 5.7 Mbp, with the MHC split into >2 contigs. The PacBio

assembly is less contiguous, despite a higher average read length and simplified haploid genome.

In addition to assembling the MHC into a single contig, the ultra-long MinION reads enabled the contiguous MHC to be haplotype phased. Due to the limited depth of nanopore reads, heterozygous SNPs were called using Illumina data and then phased using the ultra-long nanopore reads to generate two pseudo-haplotypes, from which MHC typing was performed using the approach of Dilthey *et al.*⁴⁰ (Fig. 5a). Some gaps were introduced during haplotig (contigs with the same haplotype) assembly, owing to low phased-read coverage—for example, *HLA-DRB3* was left unassembled on haplotype A—but apart from one *HLA-DRB1* allele, sample HLA types were

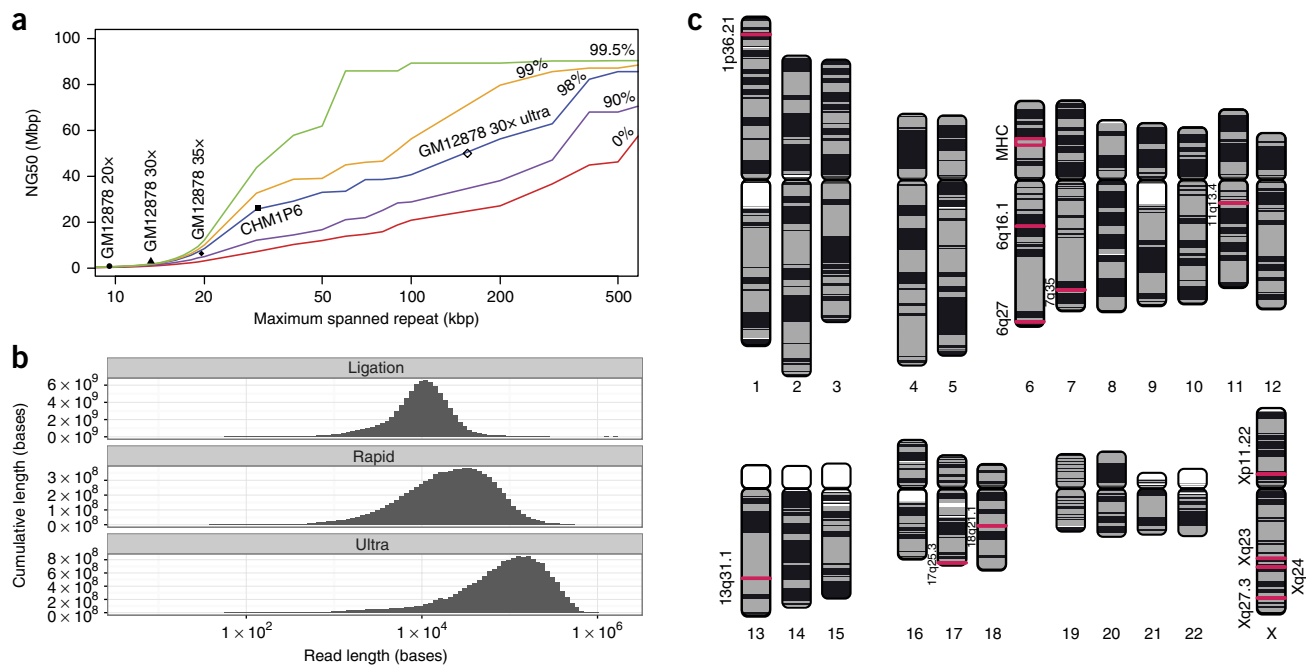


Figure 4 Repeat modeling and assembly. **(a)** A model of expected NG50 contig size when correctly resolving human repeats of a given length and identity. The y axis shows the expected NG50 contig size when repeats of a certain length (x axis) or sequence identity (colored lines) can be consistently resolved. Nanopore assembly contiguity (GM12878 20 \times , 30 \times , 35 \times) is currently limited by low coverage of long reads and a high error rate, making repeat resolution difficult. These assemblies approximately follow the predicted assembly contiguity. The *projected* assembly contiguity using 30 \times of ultra-long reads (GM12878 30 \times ultra) exceeds 30 Mbp. A recent assembly of 65 \times PacBio P6 data with an NG50 of 26 Mbp is shown for comparison (CHM1 P6). **(b)** Yield by read length (\log_{10}) for ligation, rapid and ultra-long rapid library preparations. **(c)** Chromosomes plot illustrating the contiguity of the nanopore assembly boosted with ultra-long reads. Contig and alignment boundaries, not cytogenetic bands, are represented by a color switch, so regions of continuous color indicate regions of contiguous sequence. White areas indicate unmapped sequence, usually caused by N's in the reference genome. Regions of interest, including the 12 50+ kb gaps in GRCh38 closed by our assembly as well as the MHC (16 Mbp), are outlined in red.

recovered almost perfectly with an edit distance of 0–1 for true allele versus called allele (**Supplementary Table 9**). Analysis of parental (GM12891, GM12892) HLA types confirmed the absence of switch errors between the classical HLA typing genes. To our knowledge, this is the first time the MHC has been assembled and phased over its full length in a diploid human genome.

Already published single-molecule human genome assemblies contain multiple contigs that span the MHC^{5,41,42} and phasing has not been attempted. Instead, MHC surveys have focused on homozygous cell lines⁴³.

Ultra-long reads close gaps in the human reference genome

Large (>50 kb) bridged scaffold gaps remain unresolved in the reference human genome assembly (GRCh38). These breaks in the assembly span tandem repeats and/or long tracts of segmental duplications⁴⁴. Using sequence from our *de novo*-assembled contigs, we were able to close 12 gaps, each of which was more than 50 kb in the reference genome. We then looked for individual ultra-long reads that spanned gaps, and matched the sequence closure for each region as predicted by the assembly (**Supplementary Table 10**).

The gap closures enabled us to identify 83,980 bp of previously unknown euchromatic sequence. For example, an unresolved 50-kbp scaffold gap on Xq24 marks the site of a human-specific tandem repeat that contains a cancer/testis gene family, known as CT47 (refs. 45,46). This entire region is spanned by a single contig in our final assembly (tig00002632). Inspection of this contig using hidden Markov model (HMM) profile modeling of an individual repeat unit containing the *CT47A11* gene (GRCh38 chrX:120932333–120938697) suggests that there is an array of eight tandem copies of the CT47 repeat (**Fig. 5b**). In support of this finding, we identified three ultra-long reads that together traversed the entire tandem array (**Fig. 5b**); two reads provide evidence for an array of eight repeat copies and one read supports six copies, suggesting heterozygosity.

Telomere repeat lengths

FISH (fluorescent *in situ* hybridization) estimates and direct cloning of telomeric DNA suggests that telomere repeats (TTAGGG) extend for multiple kilobases at the ends of each chromosome^{47,48}. Using HMM profile modeling of the published telomere tract of repeats (M19947.1), we identified 140 ultra-long reads that contained the TTAGGG tandem repeat (**Supplementary Table 11**). Sequences next to human telomeres are enriched in intra- and interchromosomal segmental duplications, which makes it difficult to map ultra-long reads directly to the chromosome assemblies. However, we were able to map 17/140 ultra-long reads to specific chromosome subtelomeric regions. We analyzed the mapped regions by identifying the junction or the start of the telomeric array on 17 ultra-long reads, and annotating all TTAGGG-repeat sequences to the end of the read to estimate telomeric repeat length. For example, two reads that only mapped to chromosome 21q indicate that there are 9,108 bp of telomeric repeats. Overall, we found evidence for telomeric arrays that span 2–11 kb within 14 subtelomeric regions for GM12878 (**Fig. 5c,d** and **Supplementary Table 11**).

DISCUSSION

We report sequencing and assembly of a human genome with 99.88% accuracy and an NG50 of 6.4 Mb using unamplified DNA and nanopore reads followed by short-read consensus improvement. At 30× coverage we have produced the most contiguous assembly of a human genome to date, using only a single sequencing technology

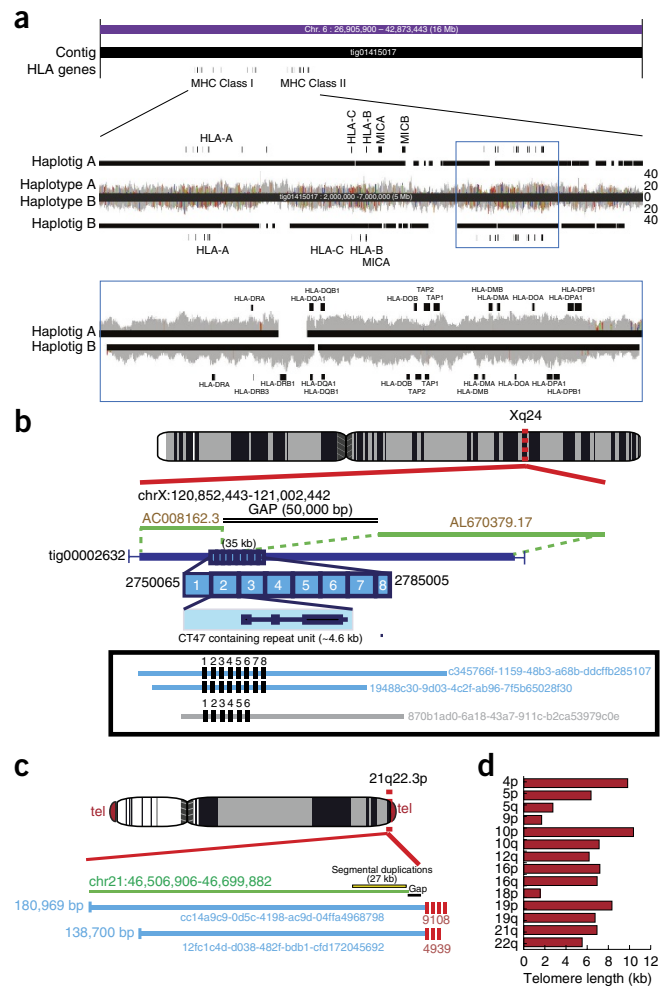


Figure 5 Ultra-long reads, assembly, and telomeres. **(a)** A 16-Mbp ultra-long read contig and associated haplotigs are shown spanning the full MHC region. MHC Class I and II regions are annotated along with various HLA genes. Below this contig, the MHC region is enlarged, showing haplotype A and B coverage tracks for the phased nanopore reads. Nanopore reads were aligned back to the polished Canu contig, with colored lines indicating a high fraction of single-nucleotide discrepancies in the read pileups (as displayed by the IGV³¹ browser). The many disagreements indicate the contig is a mosaic of both haplotypes. The haplotig A and B tracks show the result of assembling each haplotype read set independently. Below this, the MHC class II region is enlarged, with haplotype A and B raw reads aligned to their corresponding, unpolished haplotigs. The few consensus disagreements between raw reads and haplotigs indicate successful partitioning of the reads into haplotypes. **(b)** An unresolved, 50-kb bridged scaffold gap on Xq24 remains in the GRCh38 assembly (adjacent to scaffolds AC008162.3 and AL670379.17, shown in green). This gap spans a ~4.6-kb tandem repeat containing cancer/testis gene family 47 (CT47). This gap is closed by assembly (contig: tig00002632) and has eight tandem copies of the repeat, validated by alignment of 100 kb+ ultra-long reads also containing eight copies of the repeat (light blue with read name identifiers). One read has only six repeats, suggesting the tandem repeated units are variable between homologous chromosomes. **(c)** Ultra-long reads can predict telomere length. Two 100 kb+ reads that map to the subtelomeric region of the chromosome 21 q-arm, each containing 4.9–9.1 kb of the telomeric (TTAGGG) repeat. **(d)** Telomere length estimates showing variable lengths between non-homologous chromosomes.

and the Canu assembler²³. Consistent with the view that the underlying ionic raw current contains additional information, signal-based

polishing¹⁴ improved the assembly accuracy to 99.44%. Finally, we report that combining signal-based polishing and short-read (Illumina) correction²⁶ gave an assembly accuracy of 99.96%, which is similar to metrics for other mammalian genomes⁹.

Here we report that read lengths produced by the MinION nanopore sequencer were dependent on the input fragment length. We found that careful preparation of DNA in solution using classical extraction and purification methods can yield extremely long reads. The longest read lengths were achieved using the transposase-based rapid library kit in conjunction with methods of DNA extraction designed to mitigate shearing. We produced 5× coverage with ultra-long reads, and used this data set to augment our initial assembly. The final 35× coverage assembly has an NG50 of 6.4 Mb. Based on modeling we predict that 30× of ultra-long reads alone would result in an assembly with a contig NG50 in excess of 40 Mb, approaching the contiguity of the current human reference (Fig. 4c). We posit that there may be no intrinsic read-length limit for pore-based sequencers, other than from physical forces that lead to DNA fragmentation in solution. Therefore, there is scope to further improve the read-length results obtained here, perhaps through solid-phase DNA extraction and library preparation techniques, such as agar encasement.

The increased single-molecule read length that we report here, obtained using a MinION nanopore sequencer, enabled us to analyze regions of the human genome that were previously intractable with state-of-the-art sequencing methods. For example, we were able to phase megabase regions of the human genome in single contigs, to more accurately estimate telomere lengths, and to resolve complex repeat regions. Phasing of 4- to 5-Mb scaffolds through the MHC has recently been reported using a combination of sequencing and genealogical data⁴⁹. However, the resulting assemblies contained multiple gaps of unknown sequences. We phased the entire MHC, and reconstructed both alleles. Development of tools to automate phasing from nanopore assemblies is now needed.

We also wrote custom software/algorithms (poredb) to track the large number of reads, store each read as an individual file, and enable use of cloud-based pipelines for our analyses.

Our proof-of-concept demonstration of human genome sequencing using a MinION nanopore sequencer reveals the potential of this approach, but identifies specific challenges for future projects. Improvements in real-time base-calling are needed to simplify the workflow. More compact and convenient formats for storing raw and base-called data are urgently required, ideally employing a standardized, streaming compatible serialization format such as BAM/CRAM.

With ultra-long reads we found the longest reads exceeded CIGAR string limitations in the BAM format, necessitating the use of SAM or CRAM (<https://github.com/samtools/hts-specs/issues/40>). And, we were unable to complete an alignment of the ultra-long reads using BWA-MEM, and needed to adopt other algorithms, including GraphMap and NGM-LR, to align the reads. This required large amounts of compute time and RAM^{37,38,50}. Availability of our data set has spurred the development of Minimap2 (ref. 39), and we recommend this long-read aligner for use in aligning ultra-long reads on a standard desktop computer.

Nanopore genotyping accuracy currently lags behind short-read sequencing instruments, due to a limited ability to discriminate between heterozygous and homozygous alleles, which arose from error rate and the depth of coverage in our sequencing data. We found that >99% of SNP calls were correct at homozygous reference sites, dropping to 91.4% at heterozygous and homozygous non-reference sites. Similarly, Nanopore and Illumina SV genotypes agreed at 81% of

heterozygous and 90% of homozygous sites. These results highlight a need for structural variant genotyping tools for long, single-molecule sequencing reads. Using 1D² chemistry (which sequences template and complement strands of the same molecule) or modeling nanopore ionic raw current, perhaps by incorporating training data from modified DNA, could potentially produce increased read accuracy. A complementary approach would be to increase coverage.

In summary, we provide evidence that a portable, biological nanopore sequencer could be used to sequence, assemble, and provisionally analyze structural variants and detect epigenetic marks, in point-of-care human genomics applications in the future.

METHODS

Methods, including statements of data availability and any associated accession codes and references, are available in the [online version of the paper](#).

Note: Any Supplementary Information and Source Data files are available in the online version of the paper.

ACKNOWLEDGMENTS

We acknowledge the support of Oxford Nanopore Technologies staff in generating this data set, in particular R. Dokos, O. Hartwell, J. Pugh, and C. Brown. We thank M. Akeson for his support and insight. We thank R. Poplawski and S. Thompson for technical assistance with configuring and using cloud-based file systems with millions of files on CLIMB. We thank W. Timp and R. Workman for generating the R9.4 methylation training data for nanopore. We thank T. Allers for assistance with PFGE. We thank A. Pizarro at Amazon Web Services for hosting the human genome data set as an Amazon Web Services Open Data set. This study utilized the computational resources of the NIH HPC Biowulf cluster (<https://hpc.nih.gov>). This study was partially supported by the UK Antimicrobial Resistance Cross Council Initiative (JOG: MR/N013956/1), Rosetrees Trust (JOG: A749), the Intramural Research Program of the National Human Genome Research Institute, National Institutes of Health (S.K., A.D., A.R., A.M.P.), the BBSRC (M.L.: BB/N017099/1 and BB/M020061/1), the Canadian Institutes of Health Research (J.R.T., T.P.S.: #10677), Brain Canada Multi-Investigator Research Initiative Grant with Genome British Columbia the Michael Smith Foundation for Health Research, and the Koerner Foundation (J.R.T., T.P.S.), the Canada Research Chair in Biotechnology and Genomics-Neurobiology (T.P.S.), the Ontario Institute for Cancer Research and the Government of Canada (J.T.S.: OGI-129), the US National Cancer Institute (A.R.Q., B.S.P., T.A.S.: NIH U24CA209999), the Wellcome Trust (A.D.B.: 102732/Z/13/Z, M.L.: 204843/Z/16/Z), Cancer Research UK (A.D.B.: A23923), the MRC (A.D.B.: MR/M016587/1), the MRC Fellowship in Microbial Bioinformatics as part of CLIMB (N.J.L.) and the NIHR Surgical Reconstruction and Microbiology Research Centre (SRMRC).

AUTHOR CONTRIBUTIONS

N.J.L., M.L., J.T.S., and J.R.T. conceived the study, J.Q. developed the long read protocol, A.D.B., M.J., M.L., H.M., S.M., T.N., J.O'G., J.Q., H.R., J.R.T. and L.T. prepared materials and/or performed sequencing, A.T.D., I.T.F., M.J., S.K., N.J.L., M.L., K.H.M., H.E.O., B.P., B.S.P., A.M.P., A.R.Q., A.C.R., A.R., T.A.S., J.T.S. and J.R.T. performed bioinformatics analysis and wrote or modified software, I.T.F., M.J., S.K., N.J.L., M.L., K.H.M., J.O'G., H.E.O., B.P., A.M.P., J.Q., A.R.Q., A.C.R., T.A.S., J.T.S., T.P.S., and J.R.T. wrote and edited the manuscript. All authors approved the manuscript and provided strategic oversight for the work.

COMPETING FINANCIAL INTERESTS

The authors declare competing financial interests: details are available in the [online version of the paper](#).

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



This work is licensed under a Creative Commons Attribution 4.0 International licence. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons licence, users will need to obtain permission from the licence holder to reproduce the material. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

1. Wheeler, D.A. *et al.* The complete genome of an individual by massively parallel DNA sequencing. *Nature* **452**, 872–876 (2008).
2. Bentley, D.R. *et al.* Accurate whole human genome sequencing using reversible terminator chemistry. *Nature* **456**, 53–59 (2008).
3. Pushkarev, D., Neff, N.F. & Quake, S.R. Single-molecule sequencing of an individual human genome. *Nat. Biotechnol.* **27**, 847–850 (2009).
4. Rothberg, J.M. *et al.* An integrated semiconductor device enabling non-optical genome sequencing. *Nature* **475**, 348–352 (2011).
5. Pendleton, M. *et al.* Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nat. Methods* **12**, 780–786 (2015).
6. Warburton, P.E. *et al.* Analysis of the largest tandemly repeated DNA families in the human genome. *BMC Genomics* **9**, 533 (2008).
7. Wevrick, R. & Willard, H.F. Long-range organization of tandem arrays of alpha satellite DNA at the centromeres of human chromosomes: high-frequency array-length polymorphism and meiotic stability. *Proc. Natl. Acad. Sci. USA* **86**, 9394–9398 (1989).
8. 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).
9. Gordon, D. *et al.* Long-read sequence assembly of the gorilla genome. *Science* **352**, aae0344 (2016).
10. Chaisson, M.J.P., Wilson, R.K. & Eichler, E.E. Genetic variation and the de novo assembly of human genomes. *Nat. Rev. Genet.* **16**, 627–640 (2015).
11. Jain, M., Olsen, H.E., Paten, B. & Akeson, M. The Oxford Nanopore MinION: delivery of nanopore sequencing to the genomics community. *Genome Biol.* **17**, 239 (2016).
12. Quick, J. *et al.* Real-time, portable genome sequencing for Ebola surveillance. *Nature* **530**, 228–232 (2016).
13. Quick, J. *et al.* Rapid draft sequencing and real-time nanopore sequencing in a hospital outbreak of Salmonella. *Genome Biol.* **16**, 114 (2015).
14. 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).
15. Istace, B. *et al.* de novo assembly and population genomic survey of natural yeast isolates with the Oxford Nanopore MinION sequencer. *Gigascience* **6**, 1–13 (2017).
16. Datema, E. *et al.* The megabase-sized fungal genome of *Rhizoctonia solani* assembled from nanopore reads only. Preprint at <https://www.biorxiv.org/content/early/2016/11/01/084772> (2016).
17. Tyson, J.R. *et al.* Whole genome sequencing and assembly of a *Caenorhabditis elegans* genome with complex genomic rearrangements using the MinION sequencing device. Preprint at <https://www.biorxiv.org/content/early/2017/01/30/099143> (2017).
18. Quick, J. *et al.* Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly from clinical samples. Preprint at <https://www.biorxiv.org/content/early/2017/01/09/098913> (2017).
19. Jansen, H.J. *et al.* Rapid de novo assembly of the European eel genome from nanopore sequencing reads. Preprint at <https://www.biorxiv.org/content/early/2017/01/20/101907> (2017).
20. Zook, J.M. *et al.* Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
21. Jain, M. *et al.* Improved data analysis for the MinION nanopore sequencer. *Nat. Methods* **12**, 351–356 (2015).
22. Durbin, R., Eddy, S.R., Krogh, A. & Mitchison, G. *Biological Sequence Analysis* (Cambridge University Press, 1998).
23. Eberle, M.A. *et al.* A reference data set of 5.4 million phased human variants validated by genetic inheritance from sequencing a three-generation 17-member pedigree. *Genome Res.* **27**, 157–164 (2017).
24. 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).
25. Schneider, V.A. *et al.* Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res.* **27**, 849–864 (2017).
26. Walker, B.J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
27. Bickhart, D.M. *et al.* Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome. *Nat. Genet.* **49**, 643–650 (2017).
28. Chiang, C. *et al.* SpeedSeq: ultra-fast personal genome analysis and interpretation. *Nat. Methods* **12**, 966–968 (2015).
29. Layer, R.M., Chiang, C., Quinlan, A.R. & Hall, I.M. LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol.* **15**, R84 (2014).
30. Zook, J.M. *et al.* Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. *Nat. Biotechnol.* **32**, 246–251 (2014).
31. Robinson, J.T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
32. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
33. Wescoe, Z.L., Schreiber, J. & Akeson, M. Nanopores discriminate among five C5-cytosine variants in DNA. *J. Am. Chem. Soc.* **136**, 16582–16587 (2014).
34. Laszlo, A.H. *et al.* Detection and mapping of 5-methylcytosine and 5-hydroxymethylcytosine with nanopore MspA. *Proc. Natl. Acad. Sci. USA* **110**, 18904–18909 (2013).
35. Rand, A.C. *et al.* Mapping DNA methylation with high-throughput nanopore sequencing. *Nat. Methods* **14**, 411–413 (2017).
36. Simpson, J.T. *et al.* Detecting DNA cytosine methylation using nanopore sequencing. *Nat. Methods* **14**, 407–410 (2017).
37. Sović, I. *et al.* Fast and sensitive mapping of nanopore sequencing reads with GraphMap. *Nat. Commun.* **7**, 11307 (2016).
38. Sedlazeck, F.J. *et al.* Accurate detection of complex structural variations using single molecule sequencing. Preprint at <https://www.biorxiv.org/content/early/2017/07/28/169557> (2017).
39. Li, H. Minimap2: fast pairwise alignment for long DNA sequences. Preprint at <https://arxiv.org/abs/1708.01492> (2017).
40. Dilthey, A.T. *et al.* High-accuracy HLA type inference from whole-genome sequencing data using population reference graphs. *PLOS Comput. Biol.* **12**, e1005151 (2016).
41. Shi, L. *et al.* Long-read sequencing and de novo assembly of a Chinese genome. *Nat. Commun.* **7**, 12065 (2016).
42. Seo, J.-S. *et al.* De novo assembly and phasing of a Korean human genome. *Nature* **538**, 243–247 (2016).
43. Norman, P.J. *et al.* Sequences of 95 human MHC haplotypes reveal extreme coding variation in genes other than highly polymorphic HLA class I and II. *Genome Res.* **27**, 813–823 (2017).
44. Bovee, D. *et al.* Closing gaps in the human genome with fosmid resources generated from multiple individuals. *Nat. Genet.* **40**, 96–101 (2008).
45. Chen, Y.-T. *et al.* Identification of a new cancer/testis gene family, CT47, among expressed multicopy genes on the human X chromosome. *Genes Chromosom. Cancer* **45**, 392–400 (2006).
46. Jain, M. *et al.* Improved data analysis for the MinION nanopore sequencer. *Nat. Methods* **12**, 351–356 (2015).
47. Moyzis, R.K. *et al.* A highly conserved repetitive DNA sequence, (TTAGGG)_n, present at the telomeres of human chromosomes. *Proc. Natl. Acad. Sci. USA* **85**, 6622–6626 (1988).
48. Kimura, M. *et al.* Measurement of telomere length by the Southern blot analysis of terminal restriction fragment lengths. *Nat. Protoc.* **5**, 1596–1607 (2010).
49. Maretty, L. *et al.* Sequencing and de novo assembly of 150 genomes from Denmark as a population reference. *Nature* **548**, 87–91 (2017).
50. Jain, C., Dilthey, A., Koren, S., Aluru, S. & Phillippy, A.M. in *Proc. 21st Annual International Conference, RECOMB 2017* (ed., Sahinalp, S.) 66–81 (Springer, 2017).

ONLINE METHODS

Human DNA. Human genomic DNA from the GM12878 human cell line (CEPH/Utah pedigree) was either purchased from Coriell as DNA (cat. no. NA12878) or extracted from the cultured cell line also purchased from Coriell (cat. no. GM12878). Cell culture was performed using Epstein–Barr virus (EBV)-transformed B lymphocyte culture from the GM12878 cell line in RPMI-1640 media with 2 mM L-glutamine and 15% FBS at 37 °C.

QIAGEN DNA extraction. DNA was extracted from cells using the QIAamp DNA mini kit (Qiagen). 5×10^6 cells were spun at 300g for 5 min to pellet. The cells were resuspended in 200 μ l PBS and DNA was extracted according to the manufacturer's instructions. DNA quality was assessed by running 1 μ l on a genomic ScreenTape on the TapeStation 2200 (Agilent) to ensure a DNA Integrity Number (DIN) >7 (value for NA12878 was 9.3). Concentration of DNA was assessed using the dsDNA HS assay on a Qubit fluorometer (Thermo Fisher).

Library preparation (SQK-LSK108 1D ligation genomic DNA). 1.5–2.5 μ g human genomic DNA was sheared in a Covaris g-TUBE centrifuged at 5,000–6,000 r.p.m. in an Eppendorf 5424 (or equivalent) centrifuge for 2×1 min, inverting the tube between centrifugation steps.

DNA repair (NEBNext FFPE DNA Repair Mix, NEB M6630) was performed on purchased DNA but not on freshly extracted DNA. 8.5 μ l nuclease-free water (NFW), 6.5 μ l FFPE Repair Buffer and 2 μ l FFPE DNA Repair Mix were added to the 46 μ l sheared DNA. The mixture was incubated for 15 min at 20 °C, cleaned up using a 0.4 \times volume of AMPure XP beads (62 μ l), incubated at room temperature with gentle mixing for 5 min, washed twice with 200 μ l fresh 70% ethanol, pellet allowed to dry for 2 min, and DNA eluted in 46 μ l NFW or EB (10 mM Tris pH 8.0). A 1 μ l aliquot was quantified by fluorometry (Qubit) to ensure ≥ 1 μ g DNA was retained.

End repair and dA-tailing (NEBNext Ultra II End-Repair/dA-tailing Module) was then performed by adding 7 μ l Ultra II End-Prep buffer, 3 μ l Ultra II End-Prep enzyme mix, and 5 μ l NFW. The mixture was incubated at 20 °C for 10 min and 65 °C for 10 min. A 1 \times volume (60 μ l) AMPure XP clean-up was performed and the DNA was eluted in 31 μ l NFW. A 1- μ l aliquot was quantified by fluorometry (Qubit) to ensure ≥ 700 ng DNA was retained.

Ligation was then performed by adding 20 μ l Adaptor Mix (SQK-LSK108 Ligation Sequencing Kit 1D, Oxford Nanopore Technologies (ONT)) and 50 μ l NEB Blunt/TA Master Mix (NEB, cat. no. M0367) to the 30 μ l dA-tailed DNA, mixing gently and incubating at room temperature for 10 min.

The adaptor-ligated DNA was cleaned up by adding a 0.4 \times volume (40 μ l) of AMPure XP beads, incubating for 5 min at room temperature and resuspending the pellet twice in 140 μ l ABB (SQK-LSK108). The purified-ligated DNA was resuspended by adding 25 μ l ELB (SQK-LSK108) and resuspending the beads, incubating at room temperature for 10 min, pelleting the beads again, and transferring the supernatant (pre-sequencing mix or PSM) to a new tube. A 1- μ l aliquot was quantified by fluorometry (Qubit) to ensure ≥ 500 ng DNA was retained.

Sambrook and Russell DNA extraction. This protocol was modified from Chapter 6 protocol 1 of Sambrook and Russell⁵¹. 5×10^7 cells were spun at 4500g for 10 min to pellet. The cells were resuspended by pipette mixing in 100 μ l PBS. 10 ml TLB was added (10 mM Tris-Cl pH 8.0, 25 mM EDTA pH 8.0, 0.5% (w/v) SDS, 20 μ g/ml Qiagen RNase A), vortexed at full speed for 5 s and incubated at 37 °C for 1 h. 50 μ l Proteinase K (Qiagen) was added and mixed by slow inversion ten times followed by 3 h at 50 °C with gentle mixing every 1 h. The lysate was phenol-purified using 10 ml buffer saturated phenol using phase-lock gel falcon tubes, followed by phenol:chloroform (1:1). The DNA was precipitated by the addition of 4 ml 5 M ammonium acetate and 30 ml ice-cold ethanol. DNA was recovered with a glass hook followed by washing twice in 70% ethanol. After spinning down at 10,000g, ethanol was removed followed by 10 min drying at 40 °C. 150 μ l EB (Elution Buffer) was added to the DNA and left at 4 °C overnight to resuspend.

Library preparation (SQK-RAD002 genomic DNA). To obtain ultra-long reads, the standard Rapid Adapters (RAD002) protocol (SQK-RAD002 Rapid Sequencing Kit, ONT) for genomic DNA was modified as follows.

16 μ l of DNA from the Sambrook extraction at approximately 1 μ g/ μ l, manipulated with a cut-off P20 pipette tip, was placed in a 0.2 ml PCR tube, with 1 μ l removed to confirm quantification value. 5 μ l FRM was added and mixed slowly ten times by gentle pipetting with a cut-off pipette tip moving only 12 μ l. After mixing, the sample was incubated at 30 °C for 1 min followed by 75 °C for 1 min on a thermocycler. After this, 1 μ l RAD and 1 μ l Blunt/TA ligase was added with slow mixing by pipetting using a cut-off tip moving only 14 μ l ten times. The library was then incubated at room temperature for 30 min to allow ligation of RAD. To load the library, 25.5 μ l RBF (Running Buffer with Fuel mix) was mixed with 27.5 μ l NFW, and this was added to the library. Using a P100 cut-off tip set to 75 μ l, this library was mixed by pipetting slowly five times. This extremely viscous sample was loaded onto the “spot on” port and entered the flow cell by capillary action. The standard loading beads were omitted from this protocol owing to excessive clumping when mixed with the viscous library.

MinION sequencing. MinION sequencing was performed as per manufacturer's guidelines using R9/R9.4 flow cells (FLO-MIN105/FLO-MIN106, ONT). MinION sequencing was controlled using Oxford Nanopore Technologies MinKNOW software. The specific versions of the software used varied from run to run but can be determined by inspection of fast5 files from the data set. Reads from all sites were copied off to a volume mounted on a CLIMB virtual server (<http://www.climb.ac.uk>) where metadata was extracted using poredb (<https://github.com/nickloman/poredb>) and base-calling performed using Metrichor (predominantly workflow ID 1200, although previous versions were used early on in the project) (<http://www.metrichor.com>). We note that base-calling in Metrichor has now been superseded by Albacore and is no longer available. Scrapie (<https://github.com/nanoporetech/scrapie>) was used for the chr20 comparisons using reads previously identified as being from this chromosome after mapping the Metrichor reads. Albacore 0.8.4 (available from the Oxford Nanopore Technologies user community) was used for the ultra-long read set, as this software became the recommended base-caller for nanopore reads in March 2017. Given the rapid development of upgrades to base-caller software we expect to periodically re-base-call these data and make the latest results available to the community through the Amazon Open Data site.

Modified MinION running scripts. In a number of instances, MinION sequencing control was shifted to customized MinKNOW scripts. These scripts provided enhanced pore utilization/data yields during sequencing, and operated by monitoring and adjusting flow cell bias-voltage (–180 mV to –250 mV), and used an event-yield-dependent (70% of initial hour in each segment) initiation of active pore channel assignment via remuxing (reselection of ideal pores for sequencing from each group of four wells available around each channel on the flowcell). More detailed information on these scripts can be found on the Oxford Nanopore Technologies user community. In addition, a patch for all files required to modify MinION running scripts compatible with MinKNOW 1.3.23 only is available (**Supplementary Code 1**).

Live run monitoring. To assist in choosing when to switch from a standard run script to a modified run protocol, a subset of runs was monitored with the assistance of the minControl tool, an alpha component of the minoTour suite of MinION run and analysis tools (<https://github.com/minoTour/minoTour>). minControl collects metrics about a run directly from the grouper software, which runs behind the standard ONT MinKNOW interface. minControl provides a historical log of yield measured in events from a flow cell enabling estimations of yield and the decay rate associated with loss of sequencing pores over time. MinKNOW yield is currently measured in events and is scaled by approximately 1.7 to estimate yield in bases.

Assembly. All “NG” statistics were computed using a genome size of 3,098,794,149 bp (3.1 Gbp), the size of GRCh38 excluding alt sites. Canu v1.4 (+11 commits) r8006 (4a7090bd17c914f5c21bacbef4add163e492d54) was used to assemble the initial 20-fold coverage data set:

```
canu -p asm -d asm genomeSize=3.1g gridOptionsJobName=na12878nano "gridOptions==time 72:00:00-partition norm" -nanopore-raw rel2*.fastq.gz corMinCoverage=0 corMaxEvidenceErate=0.22 errorRate=0.045
```

These are the suggested low-coverage parameters from the Canu documentation, but with a decreased maximum evidence error rate. This specific parameter was decreased to reduced memory requirements after it was determined that the MinHash overlapping algorithm was underestimating error rates owing to systematic error in the reads. Counterintuitively, this systematic error makes two reads look more similar than they are, because they share more k -mers than expected under a random model. Manually decreasing the maximum overlap error rate threshold adjusted for this bias. The assembly took 40K CPU hours (25K to correct and 15K to assemble). This is about two-fold slower than a comparable PacBio data set, mostly because of the higher noise and errors in the nanopore reads.

The same version of Canu was also used to assemble the 30-fold data set:

```
canu -p asm -d asm genomeSize=3.1g gridOptionsJobName=na12878nano "gridOptions=-time 72:00:00-partition norm" -nanopore-raw rel3*.fastq.gz corMinCoverage=0 corMaxEvidenceErate=0.22 errorRate=0.045 "corMhapOptions=-threshold 0.8-num-hashes 512-ordered-sketch-size 1000-ordered-kmer-size 14"
```

For this larger data set, overlapping was again tweaked by reducing the number of hashes used and increasing the minimum overlap identity threshold. This has the effect of lowering sensitivity to further compensate for the bias in the input reads. This assembly required 62K CPU hours (29K to correct, 33K to assemble) and a peak of 120 Gbp of memory, which is about fourfold slower than a comparable PacBio data set. The assembly ran on a cluster comprised of a mix of 48-thread dual-socket Intel E5-2680 v3 @ 2.50GHz CPUs with 128 Gbp of memory and 8-thread dual-socket Intel CPU E5-2698 v4 @ 2.20GHz CPUs with 1,024 Gbp of memory.

The combined data set incorporating an additional 5× coverage of ultra-long reads was assembled with an updated version of Canu v1.4 (+125 commits) r8120:

```
canu -p asm -d asm genomeSize=3.1g gridOptionsJobName=na12878nano "gridOptions=-time 72:00:00-partition norm" -nanopore-raw rel3*.fastq.gz -nanopore-raw rel4*.fastq.gz "corMhapOptions=-threshold 0.8-num-hashes 512-ordered-sketch-size 1000-ordered-kmer-size 14" batOptions="-dg 3 -db 3 -dr 1 -el 2000 -nofilter suspicious-lopsided"
```

This assembly required 151K CPU hours (15K to correct, 86K to trim, and 50K to assemble) and a peak of 112 Gbp of memory. These high runtimes are a consequence of the ultra-long reads. In particular, the current Canu trimming algorithm was not designed for reads of this extreme length and high error rate after correction and the algorithms used are not optimal.

Assembly contiguity modeling. Expected assembly contiguity was modeled on repeat tracks downloaded from the UCSC genome browser (<http://hgdownload.soe.ucsc.edu/goldenPath/hg38/database/>).

For a given repeat identity (0%, 90%, 95%, 98%, 99%, and 99.5%), all repeats with a lower identity estimate (genomicSuperDups and chainSelf) were filtered and overlapping repeats were merged. Gaps in the reference were also considered as repeats. To compute the maximum repeat length likely to be spanned by a given sequence distribution, the probability of an unspanned repeat of a fixed length was estimated for all lengths between 1 and 100 kbp in steps of 1 kbp using an equation from <http://data-science-sequencing.github.io/lectures/lecture7/52-54>:

$$P(\text{at least one repeat is unbridged}) \leq \left(e^{-2c + \left(\frac{2N}{G}\right)} \right) \left(\sum_{i=1}^{L-2} a_i e^{\frac{2i}{G}} \right)$$

where G is the genome size, L is the read length, a_i is the number of repeats of length $1 \leq i \leq L-2$, N is the number of reads $\geq L$, and c is the coverage in reads $\geq L$. We used the distribution of all repeats for a_i and plotted the shortest repeat length such that $P(\text{at least one repeat is unbridged}) > 0.05$ for real sequencing length

distributions both nanopore and PacBio sequencing runs. Assemblies of the data were plotted at their predicted spanned read length on the x axis and NG50 on the y axis for comparison with the model. A 30× run of ultra-long coverage was simulated from the 5× dataset by repeating each ultra-long read six times.

Assembly validation and structural variant analysis. Assemblies were aligned using MUMmer v3.23 with parameters “-l 20 -c 500 -maxmatch” for the raw assemblies and “-l 100 -c 500 -maxmatch” for the polished assemblies. Output was processed with dnadiff to report average 1-to-1 alignment identity. The MUMmer coords file was converted to a tiling using the scripts from Berlin *et al.*⁵⁵ with the command:

```
python convertToTiling.py 10000 90 100000
```

and drawn using the coloredChromosomes package⁵⁶. Since the reference is a composite of human genomes and there are true variations between the reference and NA12878, we also computed a reference-free estimate of identity. A 30-fold subset of the Genome In a Bottle Illumina data set for NA12878 (ref. 20) was downloaded from ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NIST_NA12878_HG001_HiSeq_300x/RMNISTHS_30xdownsample.bam. Samtools fastq was used to extract fastq paired-end data for the full data set and for the reads mapping to chromosome 20. The reads were aligned to the whole genome assembly and chromosome 20 assemblies with BWA-MEM 0.7.12-r1039. BWA-MEM is a component of the BWA package and was chosen because of its speed and ubiquitous use in sequence mapping and analysis pipelines. Aside from the difficulties of mapping the ultra-long reads unique to this work, any other mapper could be used instead. Variants were identified using FreeBayes v1.0.2 (ref. 57), a widely used method originally developed for short-read sequencing but also applicable to long reads, with the command:

```
freebayes -C 2 -O -O -q 20 -z 0.10 -E 0 -X -u -p 2 -F 0.6 -b alignments.bam -v asm.bayes.vcf -f asm.fasta
```

The length of all variants was summed and the total number of bases with at least 3× coverage was summed using samtools depth. QV was computed as $-10 \log_{10} \left(\frac{\text{length of variants}}{\# \text{bases} \geq 3X \text{coverage}} \right)$ and identity was computed as

$100 * \left(- \frac{\text{length of variants}}{\# \text{bases} \geq 3X \text{coverage}} \right)$ Dotplots were generated with “mummer-

plot-fat” using the 1-to-1 filtered matches.

A previously published GM12878 PacBio assembly⁵ was aligned as above with MUMmer v3.23. The resulting alignment files were uploaded to Assemblytics⁵⁸ to identify structural variants and generate summary figures. Versus GRCh38, the PacBio assembly identified 10,747 structural variants affecting 10.84 Mbp, and reported an equal balance of insertions and deletions (2,361 vs. 2,724), with a peak at approximately 300 bp corresponding to Alu repeats (**Supplementary Fig. 5a** and **Supplementary Table 6**). The high error rate of the nanopore assembly resulted in a much larger number of identified variants (69,151) affecting 23.45 Mbp, with a strong deletion bias (3,900 insertions vs. 28,791 deletions) (**Supplementary Fig. 5b** and **Supplementary Table 6**). The Illumina-polished assembly reduced the total variants (47,073) affecting 16.24 Mbp but the deletion bias persisted (2,840 insertions vs. 20,797 deletions) (**Supplementary Fig. 5c** and **Supplementary Table 6**).

Base-call analysis. Sequences were aligned to the 1000 Genome GRCh38 reference (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/GRCh38_reference_genome/GRCh38_full_analysis_set_plus_decoy_hla.fa.sa) using BWA-MEM version 0.7.12-r1039 with the “-x ont2d” option⁵⁹. The BAM alignments were converted to PAF format⁶⁰ and CIGAR-strings parsed to convert alignments to an identity. Summary statistics for each flow cell were tabulated separately and combined. Alignment length versus identity was plotted using smoothScatter in R. Depth of coverage statistics for each flow cell were obtained from “samtools depth -a” and combined. As for the assembly statistics, a genome size of 3,098,794,149 bp was used to compute

bases covered. The mean coverage was 25.63 (63.20 s.d.). The minimum coverage was 0 and the maximum was 44,391. Excluding 0-coverage regions, the mean coverage was 27.41 (64.98 s.d.). The coverage histogram was plotted compared with randomly generated Poisson values generated with R's `rpois` function with $\bar{x} = 27.4074$.

Metrichor reads mapping to human chromosome 20 were additionally base-called with Scrapie v0.2.7. Scrapie reads composed primarily of low-complexity sequence were identified using the `sdust` program included with Minimap (commit: 17d5bd12290e0e8a48a5df5afaeaf4d171aa133)⁶⁰ with default parameters (-w 64 -t 20). The total length of the windows in a single sequence were merged and divided by read length to compute percentage of low-complexity sequence in each read. Any read for which this percentage exceeded 50% was removed from downstream analysis. Without this filtering, BWA-MEM did not complete mapping the sequences after >30 days of runtime on 16-cores. Similar filtering on the Metrichor-based reads had only a limited effect on the data set.

To measure homopolymer accuracy, we extracted pairwise read-to-reference alignments for reads spanning all homopolymers of length 2 or greater. For efficiency, at most 1,000 randomly selected instances were considered for each homopolymer length. Each homopolymer so identified is enclosed by two non-homopolymer “boundary” bases (e.g., the T and G in TAAAG). The number of match, mismatch, insertion, and deletion alignment operations between the boundary bases was tabulated for each homopolymer, and alignments not anchored at the boundary bases with match/mismatch operations were ignored. Homopolymer call length was reported as the number of inserted bases minus the number of deleted bases in the extracted alignment, quantifying the difference between expected and observed sequence length. All base callers with the exception of Scrapie failed in large homopolymer stretches (e.g., **Supplementary Fig. 3**), consistently capping homopolymers at 5 bp (the *k*-mer length of the model). Scrapie shows significant improvement, but tended to slightly overcall short homopolymers and undercall longer ones (**Fig. 2b**).

To quantify deviations from the expected 50:50 allele ratio at heterozygous sites, 25,541 homozygous and 46,098 heterozygous SNP positions on chromosome 20 were extracted from the Illumina Platinum Genomes project VCF for GM12878, requiring a minimum distance of 10 bp between SNP positions. Scrapie base calls at these positions were extracted using `samtools mpileup`. Deviation from the expected allelic ratio was defined as $d = \text{abs}(0.5 - [\text{allele A coverage}]/[\text{allele A coverage} + \text{allele B coverage}])$. Averaged over all evaluated heterozygous SNPs, $d = 0.13$ and 90% of SNPs have $d \leq 0.27$ (corresponding to approximately $\geq 25\%$ coverage on the minor allele). Results were similar when stratified by SNP type.

Assembly polishing with nanopolish. We ran the nanopolish consensus-calling algorithm¹⁴ on the chromosome 20 assemblies described above. For each assembly we sampled candidate variants from the base-called reads used to construct the contigs (using the “-alternative-basecalls” option) and input the original fast5 files (generated by the base-caller in the Metrichor computing platform) into a hidden Markov model, as these files contained the annotated events that the HMM relies on. The reads were mapped to the draft assembly using BWA-MEM with the “-x ont2d” option.

Each assembly was polished in 50,000-bp segments, and the individual segments were merged into the final consensus. The nanopolish jobs were run using default parameters except the “-fix-homopolymers” and “-min-candidate-frequency 0.01” options were applied.

Assembly annotation. Comparative Annotation Toolkit (CAT) (<https://github.com/ComparativeGenomicsToolkit/Comparative-Annotation-Toolkit/commit/c9503e7ad7718a935b10a72f75302caa5accb15e>) was run on both the polished and unpolished assemblies. CAT uses whole genome alignments to project transcripts from a high-quality reference genome to other genomes in the alignment⁶¹. The gene finding tool AUGUSTUS is used to clean up these transcript projections and a combined gene set is generated⁶².

To guide the annotation process, we obtained human RNA-seq data from SRA for a variety of tissues (**Supplementary Table 7**) and aligned them to both GRCh38 and the two assembly versions. GENCODE V24 was used as the reference

annotation. Two separate progressiveCactus⁶³ alignments were generated for each assembly version with the chimpanzee genome as an outgroup.

The frequency of frameshifting insertions or deletions (indels) in transcripts was evaluated by performing pairwise CDS (coding DNA sequence) sequence alignments using BLAT in a translated protein parameterization. Alignments were performed both on raw `transMap` output as well as on the final consensus transcripts.

Paralogous alignments of a source transcript were resolved through a heuristic combination of alignment coverage, identity, and synteny. Synteny is measured by counting how many gene projections near the current projection match the reference genome. In the case where multiple isoforms of a gene end up in different loci as the result of this process, a rescuing process is performed that chooses the highest scoring locus to place all isoforms at so that isoforms do not end up on different contigs. Through this process, a 1-1 orthology relationship is defined.

MHC analysis. The ultra-long assembly contains the MHC region between positions 2–6 Mb within a single 16-Mbp contig (tig01415017). Heterozygous sites were extracted by mapping Illumina reads to the polished assembly using BWA-MEM with default parameters. Alignments were post-processed according to the GATK 3.7 whole-genome variant calling pipeline, except for the “-T IndelRealigner” step using “-consensusDetermination-Model USE_READS”. The -T HaplotypeCaller parameter was used for variant calling. WhatsHap⁶⁴ was used to phase the Illumina variants with Nanopore reads reported to be contained in the contig by Canu. WhatsHap was modified to accept CRAM (<http://genome.cshlp.org/content/21/5/734.long>, <https://bitbucket.org/skoren/whatschap>) output since BAM files could not represent long CIGAR strings at the time of this analysis (<https://github.com/samtools/hts-specs/issues/40>). First, WhatsHap was run excluding any ultra-long sequences. This generated 18 phase blocks across the MHC. When ultra-long sequences were included the result was a single phase block comprising the entire MHC, supporting the utility of ultra-long reads in resolving haplotypes across large, complex regions in the genome. Nanopore reads were aligned back to the assembly using NGM-LR (CoNvex Gap-cost alignments for Long Reads)³⁸ and the combined VCF file used for phasing. Reads with more than one phasing marker were classified as haplotype A or B when >55% of their variants were in agreement (**Fig. 5a**). A new assembly was generated for haplotypes A and B using only reads assigned to each haplotype as well as reads marked homozygous. The assemblies were polished by Pilon 1.21 (ref. 26) using the SGE pipeline at <https://github.com/skoren/PilonGrid>. Pilon was given all reads mapping to the MHC.

Exon sequences belonging to the six classical HLA genes were extracted from the phased assembly, and HLA types called at G group resolution. These results were compared to GM12878 HLA type reference data. For the class I and II HLA genes, with the exception of one DRB1 haplotype, there was good agreement between the best-matching reference type and the alleles called from the assembly (edit distance 0–1). Detailed examination of HLA-DRB1, however, showed that one exon (exon 2) is different from all reference types in the assembly, a likely error in the assembly sequence.

GM12878 G group HLA types for HLA-A/B/C, HLA-DQA1, HLA-DQB1, and HLA-DRB1 are from ref. 65; the presence of exactly one HLA-DRB3 allele is expected due to linkage with HLA-DRB1 (DRB1*03 is associated with HLA-DRB3, and DRB1*01 has no DRB3/4/5 association).

Genotyping SNPs using nanopolish. Nanopolish was used for genotyping the subset of reads that mapped to human chromosome 20. The 1000 Genomes phase 3 variant set for GRCh38 was used as a reference and filtered to include only chromosome 20 SNPs that were not singletons (Allele Count ≥ 2). This set of SNPs was input into “nanopolish variants” in genotyping mode (“-genotype”). The genotyping method extends the variant calling framework previously described¹² to consider pairs of haplotypes, allowing it to be applied to diploid genomes (option “-ploidy 2”). To evaluate their accuracy, genotype calls were compared to the “platinum calls” generated by Illumina²³. When evaluating the correctness of a nanopore call, we required the log-likelihood ratio of a variant call (heterozygous or homozygous non-reference) to be at least 30, otherwise, we considered the site to be homozygous reference.

Estimating SV genotyping sensitivity. Previously identified high-confidence GM12878 SVs, validated with Moleclo and/or PacBio long reads, were used to determine genotyping sensitivity²⁹. Using LUMPY²⁸, we recalled SVs in the Platinum Genomes NA12878 Illumina data set (paired-end reads; European Nucleotide Archive, Run Accession ERR194147), intersected these calls with the aforementioned high confidence set, and genotyped the resulting calls using SVTyper²⁸ and the same Platinum alignments, generating a set of 2,414 high-confidence duplications and deletions with accompanying genotypes. Nanopore reads from all flow cells were mapped using BWA-MEM (bwa mem -k15 -W30 -r10 -B2 -O2 -L0), and then merged into release-specific BAM files. Merged BAM files were subsampled using Samtools (samtools view -s \$COVERAGE_FRACTION) to approximate coverage values as shown in **Figure 2a**. SVs were then genotyped in each subsampled BAM file using a modified version of SVTyper (<http://github.com/tomsasani/svtyper>). Generally, long nanopore reads are subject to higher rates of mismatches, insertions, and deletions than short Illumina reads. These features can result in ‘bleed-through’ alignments, where reads align past the true breakpoint of an SV⁶⁶. The modifications to SVTyper attempt to correct for the bleed-through phenomenon by allowing reads to align past the breakpoint, yet still support an alternate genotype. All modifications to SVTyper are documented in the source code available at the GitHub repository listed above (commit ID: d70de9c) (**Supplementary Code 2**). Nanopore- and Illumina-derived genotypes were then compared as a function of subsampled nanopore sequencing coverage.

The false-discovery rate of our SVTyper genotyping strategy was estimated by randomly permuting the genomic locations of the original SVs using BEDTools “shuffle”⁶⁷. Centromeric, telomeric, and “gap” regions (as defined by the UCSC Genome Browser) were excluded when assigning randomly selected breakpoints to each SV. The randomly shuffled SVs were then genotyped in Illumina and nanopore data in the same manner as before. It is expected that the alignments at shuffled SV intervals would almost always support a homozygous reference genotype. So, all instances in which Illumina data supported a homozygous reference genotype, yet the nanopore data called a non-homozygous reference genotype, were considered false positives. SV coordinates were shuffled and genotyped 1,000 times and the average false-discovery rate over all iterations was 6.4%.

Nanopore and PacBio genotyping sensitivity was compared to a subset of our high-confidence SV set. Because our high-confidence set includes only “DUP” and “DEL” variants, and the Genome in a Bottle (GIAB) PacBio SV VCF (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NA12878_PacBio_MtSinai/NA12878.sorted.vcf.gz) does not report “DUP” variants, we compared genotypes at deletions with genomic coordinates that shared reciprocal overlap of at least 0.5 between the GIAB VCF and our high-confidence SV VCF. We then compared nanopore genotypes (as determined by SVTyper) with the genotypes reported in the GIAB SV VCF. Importantly, the GIAB VCF was derived from a ~44× coverage data set, whereas our data set (containing data from both releases) represents only about ~32× coverage of the genome. Additionally, all nanopore data used in this analysis were aligned using BWA, while GIAB PacBio data were aligned using BLASR⁶⁸.

Scaling marginAlign and signalAlign data analysis pipelines. To handle the large data volume, the original marginAlign and signalAlign algorithms were ported to cloud infrastructures using the Toil batch system⁶⁹. Toil allows for computational resources to be scaled horizontally and vertically as a given experiment requires and enables researchers to perform their own experiments in identical conditions. All of the workflows used and the source code is freely available from <https://github.com/ArtRand/toil-signalAlign> and <https://github.com/ArtRand/toil-marginAlign>. Workflow diagrams are shown in **Supplementary Figure 10**.

Generating a controlled set of methylated control DNA samples. For signalAlign, DNA methylation control standards were obtained from Zymo Research (cat. no. D5013). The standards contain a whole-genome-amplified (WGA) DNA substrate that lacks methylation and a WGA DNA substrate that has been enzymatically treated so all CpG dinucleotides contain 5-methylcytosines. The two substrates were sequenced independently on two different flow cells using the sequencing protocol described above.

Otherwise, training for signalAlign and nanopolish was carried out as previously described^{35,36}.

5-methylcytosine detection with signalAlign. The signalAlign algorithm uses a variable order hidden Markov model combined with a hierarchical Dirichlet process (HMM-HDP) to infer base modifications in a reference sequence using the ionic current signal produced by nanopore sequencing⁷⁰. The ionic current signal is simultaneously influenced by multiple nucleotides as the strand passes through the nanopore. Correspondingly, signalAlign models each ionic current state as a nucleotide *k*-mer. The model allows a base in the reference sequence to have any of multiple methylation states (in this case 5-methyl cytosine or canonical cytosine). The model ties the probabilities of consistently methylated *k*-mers by configuring the HMM in a variable order meta-structure that allows for multiple paths over a reference *k*-mer depending on the number of methylation possibilities. To learn the ionic current distributions for methylated *k*-mers, signalAlign estimates the posterior mean density for each *k*-mer’s distribution of ionic currents using a Markov chain Monte Carlo (MCMC) algorithm given a set of *k*-mer-to-ionic current assignments. Using the full model, the posterior for each methylation status is calculated for all cytosines in CpG dinucleotides.

5-methylcytosine detection with nanopolish. Previous work describes using nanopolish to call 5-methylcytosine in a CpG context using a hidden Markov model³⁶. The output of the nanopolish calling procedure is a log-likelihood ratio, where a positive log-likelihood ratio indicates evidence for methylation. Nanopolish groups nearby CpG sites together and calls the group jointly, assigning the same methylation status to each site in the group. To allow comparison to the bisulfite data each such group was broken up into its constituent CpG sites, which all have the same methylation frequency. Percent-methylation was calculated by converting the log-likelihood ratio to a binary methylated/unmethylated call for each read, and calculating the fraction of reads classified as methylated. A filtered score was also computed by first filtering reads where the absolute value of the log-likelihood ratio was less than 2.5 to remove ambiguous reads.

Life Sciences Reporting Summary. Further information on experimental design is available in the **Life Sciences Reporting Summary**.

Data availability. Sequence data including raw signal files (FAST5), event-level data (FAST5), base-calls (FASTQ) and alignments (BAM) are available as an Amazon Web Services Open Data set for download from <https://github.com/nanopore-wgs-consortium/NA12878>. Nanopore raw signal files and the 35× assembly are additionally archived and available from the European Nucleotide Archive under accession PRJEB23027.

- Sambrook, J. & Russell, D.W. *Molecular Cloning: A Laboratory Manual* (Cold Spring Harbor Laboratory Press, 2001).
- Shomorony, I., Courtade, T. & Tse, D. in *2015 IEEE International Symposium on Information Theory (ISIT)* 919–923 (IEEE, 2015).
- Bresler, G., Bresler, M. & Tse, D. Optimal assembly for high throughput shotgun sequencing. *BMC Bioinformatics* **14** (Suppl. 5: S18) (2013).
- Ukkonen, E. Approximate string-matching with q-grams and maximal matches. *Theor. Comput. Sci.* **92**, 191–211 (1992).
- Berlin, K. *et al.* Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* **33**, 623–630 (2015).
- Böhringer, S.G., Böhringer, R., Schulte, D. & Epplen, T. Jörg T. A software package for drawing ideograms automatically. *Online J. Bioinform.* **1**, 51–60 (2002).
- Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. Preprint at <https://arxiv.org/abs/1207.3907> (2012).
- Nattestad, M. & Schatz, M.C. Assemblytics: a web analytics tool for the detection of variants from an assembly. *Bioinformatics* **32**, 3021–3023 (2016).
- Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997v2> (2013).
- Li, H. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics* **32**, 2103–2110 (2016).
- Zhu, J. *et al.* Comparative genomics search for losses of long-established genes on the human lineage. *PLOS Comput. Biol.* **3**, e247 (2007).

62. Stanke, M., Diekhans, M., Baertsch, R. & Haussler, D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics* **24**, 637–644 (2008).
63. Paten, B. *et al.* Cactus graphs for genome comparisons. *J. Comput. Biol.* **18**, 469–481 (2011).
64. Patterson, M. *et al.* in *Research in Computational Molecular Biology* (ed., Sharan, R.) 237–249 (Springer, 2014).
65. Dilthey, A., Cox, C., Iqbal, Z., Nelson, M.R. & McVean, G. Improved genome inference in the MHC using a population reference graph. *Nat. Genet.* **47**, 682–688 (2015).
66. Norris, A.L., Workman, R.E., Fan, Y., Eshleman, J.R. & Timp, W. Nanopore sequencing detects structural variants in cancer. *Cancer Biol. Ther.* **17**, 246–253 (2016).
67. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
68. Chaisson, M.J. & Tesler, G. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* **13**, 238 (2012).
69. Vivian, J. *et al.* Rapid and efficient analysis of 20,000 RNA-seq samples with Toil. Preprint at <https://www.biorxiv.org/content/early/2016/07/07/062497> (2016).
70. Rand, A.C. *et al.* Mapping DNA methylation with high-throughput nanopore sequencing. *Nat. Methods* **14**, 411–413 (2017).

Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

▶ Experimental design

1. Sample size

Describe how sample size was determined.

A single sample was sequenced and so this is not applicable.

2. Data exclusions

Describe any data exclusions.

No data were excluded from the study.

3. Replication

Describe whether the experimental findings were reliably reproduced.

The experiments were carried out at five distinct sites. As a single sample was being sequenced with data combined together, replication was not required.

4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

As a single sample was sequenced, no randomization was required in this study.

5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

No blinding was required for this study.

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

- | | |
|--------------------------|--|
| n/a | Confirmed |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The <u>exact sample size</u> (n) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A statement indicating how many times each experiment was replicated |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as an adjustment for multiple comparisons |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The test results (e.g. P values) given as exact values whenever possible and with confidence intervals noted |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A clear description of statistics including <u>central tendency</u> (e.g. median, mean) and <u>variation</u> (e.g. standard deviation, interquartile range) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Clearly defined error bars |

See the web collection on [statistics for biologists](#) for further resources and guidance.

▶ Software

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

7. Software

Describe the software used to analyze the data in this

All software used for data analysis are fully described in the materials and methods

study.

of the manuscript. All custom computer code is deposited in GitHub in appropriate repositories as described in the materials and methods. These are also linked from the main project GitHub page:

<https://github.com/nanopore-wgs-consortium/NA12878>

Software and versions used:

Oxford Nanopore Technologies:

MinKNOW 1.3.24 (ONT)

Metrichor (ONT)

Scrappie (<https://github.com/nanoporetech/scrappie> commit

2d5f7883a31152cf75ff77a060c751288f74e972) (ONT)

Albacore v 0.8.4 (ONT)

Nanopore Custom Tuning Scripts (Supplementary Code 1)

SVTyper (Supplementary Code 2 - <http://github.com/tomsasani/svtyper> commit d70de9c)

SignalAlign - <https://github.com/ArtRand/toil-signalAlign>

MarginAlign - <https://github.com/ArtRand/toil-marginAlign>

Poredb - <https://github.com/nickloman/poredb>

minControl - <https://github.com/minoTour/minoTour>

Comparative Annotation Toolkit (CAT) - <https://github.com/ComparativeGenomicsToolkit/Comparative-Annotation-Toolkit> commit c9503e7

Pilon 1.21 - <https://github.com/skoren/PilonGrid>

sdust - <https://github.com/lh3/minimap> commit

17d5bd12290e0e8a48a5df5afaeaf4d171aa133

Canu v1.4 - <https://github.com/marbl/canu> r8120 and r8006

(4a7090bd17c914f5c21bacbef4add163e492d54)

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

► Materials and reagents

Policy information about [availability of materials](#)

8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

This study uses NA12878 cell line and DNA which is supplied by Coriell and is approved for genome sequencing governed by the Coriell Institutional Review Board ("Coriell IRB") in accordance with DHHS regulations (45 CFR Part 46).

9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

No antibodies were used.

10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

This study uses NA12878 cell line and DNA which is supplied by Coriell and is approved for genome sequencing governed by the Coriell Institutional Review Board ("Coriell IRB") in accordance with DHHS regulations (45 CFR Part 46) and is not considered human subjects research.

b. Describe the method of cell line authentication used.

Purchased from validated source and sequenced. Source validates cells as described here: https://www.coriell.org/0/pdf/CC_Process_Flow.pdf

c. Report whether the cell lines were tested for mycoplasma contamination.

Cells are routinely screened for mycoplasma and mycoplasma contamination would be detectable via sequencing.

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

No commonly mis-identified cell lines were used.

► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

No animals were used in this study.

Policy information about [studies involving human research participants](#)

12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

This study uses NA12878 cell line and DNA which is supplied by Coriell and is approved for genome sequencing governed by the Coriell Institutional Review Board ("Coriell IRB") in accordance with DHHS regulations (45 CFR Part 46) and is not considered human subjects research.