

SCIENTIFIC DATA

OPEN

SUBJECT CATEGORIES

- » Next-generation sequencing
- » Genomics
- » Genome assembly algorithms
- » Standards

Data Descriptor: Extensive sequencing of seven human genomes to characterize benchmark reference materials

Justin M. Zook¹, David Catoe¹, Jennifer McDaniel¹, Lindsay Vang¹, Noah Spies^{1,2}, Arend Sidow², Ziming Weng², Yuling Liu², Christopher E. Mason³, Noah Alexander³, Elizabeth Henaff³, Alexa B.R. McIntyre³, Dhruva Chandramohan³, Feng Chen⁴, Erich Jaeger⁴, Ali Moshrefi⁴, Khoa Pham⁵, William Stedman⁵, Tiffany Liang⁵, Michael Saghbin⁵, Zeljko Dzakula⁵, Alex Hastie⁵, Han Cao⁵, Gintaras Deikus⁶, Eric Schadt⁶, Robert Sebra⁶, Ali Bashir⁶, Rebecca M. Truty⁷, Christopher C. Chang⁷, Natali Gulbahce⁷, Keyan Zhao⁸, Srinka Ghosh⁸, Fiona Hyland⁸, Yutao Fu⁸, Mark Chaisson⁹, Chunlin Xiao¹⁰, Jonathan Trow¹⁰, Stephen T. Sherry¹⁰, Alexander W. Zaranek¹¹, Madeleine Ball¹¹, Jason Bobe^{6,11}, Preston Estep^{11,12}, George M. Church^{11,12}, Patrick Marks¹³, Sofia Kyriazopoulou-Panagiotopoulou¹³, Grace X.Y. Zheng¹³, Michael Schnall-Levin¹³, Heather S. Ordonez¹³, Patrice A. Mudivarti¹³, Kristina Giorda¹³, Ying Sheng¹⁴, Karoline Bjarnesdatter Rypdal¹⁴ & Marc Salit^{1,2}

Received: 09 September 2015

Accepted: 15 March 2016

Published: 07 June 2016

The Genome in a Bottle Consortium, hosted by the National Institute of Standards and Technology (NIST) is creating reference materials and data for human genome sequencing, as well as methods for genome comparison and benchmarking. Here, we describe a large, diverse set of sequencing data for seven human genomes; five are current or candidate NIST Reference Materials. The pilot genome, NA12878, has been released as NIST RM 8398. We also describe data from two Personal Genome Project trios, one of Ashkenazim Jewish ancestry and one of Chinese ancestry. The data come from 12 technologies: BioNano Genomics, Complete Genomics paired-end and LFR, Ion Proton exome, Oxford Nanopore, Pacific Biosciences, SOLiD, 10X Genomics GemCode WGS, and Illumina exome and WGS paired-end, mate-pair, and synthetic long reads. Cell lines, DNA, and data from these individuals are publicly available. Therefore, we expect these data to be useful for revealing novel information about the human genome and improving sequencing technologies, SNP, indel, and structural variant calling, and *de novo* assembly.

¹National Institute of Standards and Technology, Gaithersburg, Maryland 20899, USA. ²Stanford University, Stanford, California 94305, USA. ³Department of Physiology and Biophysics, the Feil Family Brain and Mind Research Institute, and HRH Prince Alwaleed Bin Talal Bin Abdulaziz Alsaud Institute for Computational Biomedicine, Weill Medical College, Cornell University, New York, New York 10065, USA. ⁴Illumina Mission Bay, San Francisco, California 94158, USA. ⁵BioNano Genomics, San Diego, California 92121, USA. ⁶Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, New York 10029, USA. ⁷Complete Genomics Inc., Mountain View, California 94043, USA. ⁸Thermo Fisher Scientific, South San Francisco, California 94080, USA. ⁹Genome Sciences, University of Washington, Seattle, Washington 98105, USA. ¹⁰National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, 45 Center Drive, Bethesda, Maryland 20892, USA. ¹¹PersonalGenomes.org, Boston, Massachusetts 02115, USA. ¹²Harvard Medical School, Boston, Massachusetts 02115, USA. ¹³10X Genomics, Pleasanton, California 94566, USA. ¹⁴Department of Medical Genetics, Oslo University Hospital, Kirkeveien 166, Bygg 25, Oslo 0450, Norway. Correspondence and requests for materials should be addressed to J.Z. (email: jzook@nist.gov).

Design Type(s)	reference design • replicate design • protocol optimization design • individual genetic characteristics comparison design
Measurement Type(s)	whole genome sequencing
Technology Type(s)	DNA sequencing
Factor Type(s)	ethnic group
Sample Characteristic(s)	Homo sapiens • EBV-LCL cell

Background & Summary

Developing Reference Materials is a unique measurement science task, where significant resources can be expended to deeply characterize a small number of samples. Reference Materials act to calibrate, benchmark, or validate a measurement process. These samples are often the source of the scales on which we report our results (e.g., the Molar concentration of cholesterol), and they can be a physical realization of the SI units. Our ability to compare results between laboratories in most applications depends on Reference Materials.

There is a tradition of innovation in measurement science to characterize these high-impact samples¹. New technologies are used, rigorous experimental designs are employed, and exotic methods applied^{2,3}. In a virtuous cycle, existing methods are optimized and new methods are developed using reference materials as the benchmarks. Regulated applications depend on reference materials for quantitative, objective oversight; this opens new applications for a measurement technology, with great quality-of-life social benefit. With sequencing technologies and bioinformatics changing rapidly, whole genome reference materials and diverse data types like those presented here are a valuable resource for developing, improving, and assessing performance of these methods.

The National Institute of Standards and Technology (NIST)-hosted Genome in a Bottle Consortium is developing reference materials from well-characterized genomic DNA from 5 individuals (Fig. 1). These reference materials are the first of their kind, and will play key roles in the translation of genome sequencing to widespread adoption and as validation tools in clinical practice. We previously characterized high-confidence SNP, indel, and homozygous reference genotypes⁴, as well as large deletions and insertions⁵. We plan to use similar methods as well as new methods to characterize these genomes using the data described in this work.

The NIST Reference Material DNA has been characterized to an unprecedented degree. We have collected a large diverse set of data from 12 sequencing technologies and library preparation methods (Table 1). These data include high-depth paired-end short read whole genome sequencing (WGS) and whole exome sequencing (WES), long mate-pair WGS, pseudo long read ('read clouds') WGS, long read WGS, genome mapping, and exome sequencing. For each dataset, we describe the library preparation and sequencing methods, the currently available data records, and technical validation. We expect these data to be complementary to each other, so that we can use them to characterize a broad spectrum of phased variants of all sizes in as much of the genome as possible. We invite anyone to join in this open, public effort to characterize these genomes; thus, here, we describe the measurement methods and data as a public resource.

Methods

Participants

The pilot genome (NIST RM 8398) is an oft-used genome of Caucasian ancestry: NA12878 from the CEPH Utah Reference Collection. In addition, genomes from two family trios (both Mother-Father-Son) have been selected from the Personal Genomes Project (PGP), one of Ashkenazi Jewish (AJ) ancestry and the other of Han Chinese ancestry. These genomes are available as cells or extracted DNA from the Coriell Institute for Medical Research and are or will be available as DNA as NIST Reference Materials. The NIST Reference Materials are extracted DNA from large, homogenized batches of cells prepared specially by Coriell to control for any batch effects. The samples from PGP are consented more broadly for many applications, including commercial redistribution. There are already three commercial products from the same cell lines from which the NIST Reference Material DNA is prepared: AcroMetrix Oncology Hotspot Control from Thermo Fisher Scientific, GIAB HDx Reference Standards from Horizon Diagnostics, and cell line DNA with synthetic DNA spike-ins from SeraCare Life Sciences.

Illumina paired end WGS

Library preparation. For the AJ trio, Chinese son, and NA12878, libraries were prepared from 6 vials of the NIST Reference Material DNA for each individual. For the Chinese parents, a single library was prepared from genomic DNA from the Coriell Institute for Medical Research. For each Reference

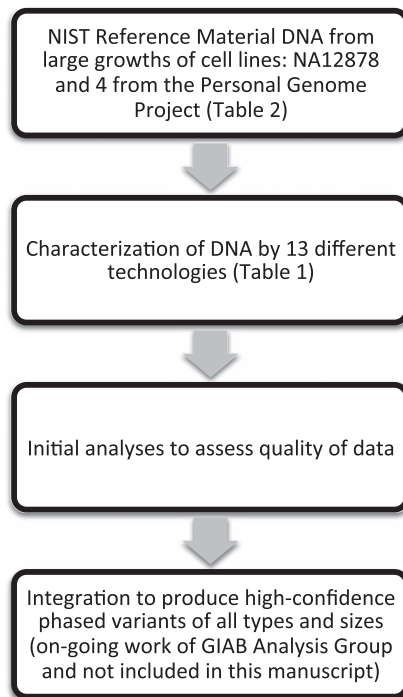


Figure 1. Overview of the study design. Data and analyses included in this manuscript are above the dotted line, and ongoing analyses of these data by the Genome in a Bottle Analysis Group are below the dotted line.

Material, 12 (or 14 for NA12878) libraries were prepared in parallel using the Illumina TruSeq (LT) DNA PCR-Free Sample Prep Kits (FC-121–3001). Two (or 3 for NA12878) libraries each were made from the first and last tubes in the lot, two libraries each were prepared from four samples pulled randomly from each quarter of the lot. This library design is intended for homogeneity analyses not presented here.

DNA concentrations were measured using a Qubit 2.0 fluorometer (Life Technologies). Genomic DNA (1.5 ug) was fragmented using a Covaris S2 focused ultrasonicator in micro TUBE AFA Fiber Pre-Slit Snap-Cap 6 × 16 mm micro tubes and the Covaris MicroTUBE holder (covaris part numbers 520045 and 500114, respectively) under the following conditions for a target insert size of 550 base pairs. Duty cycle: 10%; Intensity: 2.0; Cycles Per Burst: 200; Duration: 45 s; Mode: Frequency Sweeping; Displayed Power: 9 W; Temperature: 5.5° to 6 °C. After Fragmentation, DNA was cleaned up using illumina Sample Purification Beads. End Repair was performed in 0.2 ml PCR tubes on an MJ research PTC-200 thermal cycler. The optional end repair control was not used. Size selection was done using a 96-well 0.8 ml plate (Fisher Scientific Part # AB-0859), a magnetic stand-96 (Ambion part # AM10027) and the Illumina sample purification beads according to the 550 bp insert protocol.

Adenylation of 3' ends was done in 0.2 ml PCR tubes on an MJ Research PTC-200 thermal cycler. The optional A-Tailing control was not used. Ligation of indexed paired-end adapters was done in 0.2 ml PCR tubes using the DNA adapter tubes included in the Illumina TruSeq (LT) DNA PCR-Free Sample Prep Kit on an MJ Research PTC-200 thermal cycler. The optional ligation control was not used. The libraries were cleaned up in a 96-well 0.8 ml plate (Fisher Scientific Part # AB-0859) and a magnetic stand-96 (Ambion part # AM10027) using the Illumina sample purification beads. The final libraries were run on an Agilent 2100 Bioanalyzer HS-DNA chip to verify fragment size distribution. Final library concentration was measured via qPCR using the KAPA library quantification kit for Illumina sequencing platforms (KAPA part # KK4835). Libraries were then pooled based to the qPCR quantification data. The pool was intentionally made uneven so as to acquire greater sequence depth from the libraries made from the first and last tubes in each lot. The pools were adjusted between sequencing runs based on index balance.

For the Chinese son, DNA libraries were prepared in the same manner as they were for the Ashkenazim trio. The initial pool was made based on quantification measurements made using an Agilent 2100 Bioanalyzer, qPCR was not performed. This initial pool was sequenced on an Illumina MiSeq. The index balance obtained from the MiSeq run was used to adjust the pool for Sequencing on an Illumina HiSeq. The pool was intentionally made uneven so as to acquire greater sequence depth from the libraries made from the first and last tubes in each lot. The pools were adjusted between sequencing runs based on index balance.

Sequencing. For NA12878, the AJ Trio, and the Chinese parents, the pooled TruSeq libraries were run on an Illumina HiSeq 2500 in Rapid mode (v1) with 2 × 148 paired end reads. Pooled Libraries were

Lab	Technology	Library Type	Date Produced	Mean Paired end distance	Mean read length	Mean coverage per individual for AJ trio	Mean coverage per individual for Chinese trio
NIST	Illumina HiSeq 2500 Rapid SBS	WGS paired end PCR-free	Jul to Sep 2014	565 bp	2 × 250 (Chinese son), 2 × 148 (others)	296.83x	300x (son) 100x/parent
NIST/ Stanford	Illumina HiSeq 2500 v3/v4	WGS mate-pair	Dec 2014 to Mar 2015	6000 to 7000 bp	100	13x to 14x	14x to 17x
NIST/ Illumina	Illumina HiSeq 2500 v4	Synthetic Long Reads	Mar to Apr 2015	N/A	2 × 125	20x to 30x	20x to 30x
Oslo University Hospital	Illumina HiSeq 2500 v4	WES paired end	Oct 2015	202 bp	2 × 125	192x to 225x	181x (Son)
10X Genomics	10X GemCode WGS	WGS Linked-Reads	June to July 2015	200 to 220 bp	2 × 98	25x (son) 24x (mother) 22x (father)	None
Complete Genomics	Complete Genomics	WGS paired end	Sep 2014	388	26	101x	98x
Complete Genomics	Complete Genomics	WGS LFR	Dec 2014	278	26	100x	None
Thermo Fisher	Ion Proton	WES	Jun 2014	N/A	190	1020x	1036x
NIST	SOLiD 5500 × 1 Wildfire	WGS single end	Jan to Aug 2015	N/A	55 to 57	79x (Son)	62x (Son)
BioNano Genomics	BioNano Genomics	WG optical mapping	Feb to Dec 2014	N/A	251 kb (Chinese son), 195 kb (AJ Son), 213 kb (AJ Mother), 246 kb (AJ Father)	92x (Mother), 87x (Father), 112x (Son)	57x (Son)
NIST/Mt. Sinai	PacBio P6-C4 (90%) & P5-C3 (10%)	WGS single end	Oct 2014 to Mar 2015	N/A	10–11 kb N50	69x son; 30–32x/parent	None
Weill Cornell	Oxford Nanopore	2D reads	May 2015	N/A	5.8 kb	0.017x (son)	None

Table 1. Summary of data available for GIAB samples. Characteristics of these data for the AJ and Chinese trios are described here, and in more detail in this manuscript.

initially loaded at a concentration of 10 pM. loading concentration was adjusted accordingly on subsequent runs to balance the libraries as well as possible.

For the Chinese son, the libraries were sequenced on an Illumina HiSeq 2500 in rapid mode (v2) with 2 × 250 paired end reads. Pooled Libraries were initially loaded at a concentration based on the information from the MiSeq run. Loading concentration was adjusted accordingly on subsequent runs to optimize cluster density.

The runs were designed to get approximately 300x total coverage of each of NA12878, the AJ Trio, and the Chinese son, and 100x coverage of each of the Chinese parents.

Illumina mate-pair WGS

Library preparation. Mate Pair libraries were generated for the AJ Trio and Chinese Trio using Nextera Mate Pair Sample Preparation Kit (Illumina, Cat# FC-132-1001). Briefly, 4 µg of high molecular weight genomic DNA from the NIST Reference Materials (or from Coriell for the Chinese parents) was fragmented to about 7 kb in a 400 ml tagmentation reaction containing 12 µl of Tagment Enzyme at 55 °C for 30 min. The tagmented DNA fragments were purified with Zymo Genomic DNA Clean & Concentrator Kit (Zymo Research, Cat# D4010). The gap in the tagmented DNA was filled with a Strand Displacement Polymerase in a 200 µl strand displacement reaction at 20 °C for 30 min. DNA was then purified with AMPure XP Beads (0.5x vol, Beckman Coulter, Cat# A63880) and size-selected by 0.6% agarose gel electrophoresis in 0.5x TBE buffer. The 6–9 kb fragments were excised from gel and DNA was recovered using a Zymoclean™ Large Fragment DNA Recovery Kit (Zymo Research, Cat# D4045). Up to 600 µg of DNA was then circulated overnight at 30 °C with Circularization Ligase in a 300 µl reaction.

After overnight circularization, the uncirculated linear DNA was removed by Exonuclease digestion. Both DNA Ligase and Exonuclease were inactivated by heat treatment and the addition of Stop Ligation Buffer. Circularized DNA was then sheared to smaller sized fragments (300–1000 bp) using Covaris S2 with T6 (6 × 32 mm) glass tube (Covaris, Part# 520031 and 520042) under these conditions: Intensity of 8, Duty Cycle of 20%, Cycles Per Burst of 200, Time of 40 s, Temperature of 6–8 °C.

The sheared DNA fragments that contain the biotinylated junction adapter are mate pair fragments. These fragments were isolated by binding to Dynabeads M-280 Streptavidin Magnetic Beads (Invitrogen, Part# 112-05D) in Bead Bind Buffer. The unbiotinylated molecules in solution are unwanted genomic fragments that are removed through a series of washes. All downstream reactions were carried out on bead and beads were washed between successive reactions. The sheared DNA was first end-repaired to generate blunt ends followed by an A-Tailing reaction to add a single 'A' nucleotide to the 3' ends of the blunt fragments. Then the Illumina T-tailed indexing adapters were ligated to the A-tailed fragments.

The adapter-ligated fragments were PCR amplified [98 °C/1 min, 11 cycles of (98 °C/10 s, 60 °C/30 s, 72 °C/30 s), 72 °C/5 min, 4 °C /hold] to generate the final library. The amplified library was purified using AMPure XP Beads (0.67x vol) and eluted in Resuspension Buffer. The size distribution of the library was determined by running a sample on an Agilent Technologies 2100 Bioanalyzer. Library concentration was measured by the Qubit dsDNA HS Assay Kit (Life Technologies, Cat# Q32851).

Sequencing. Pooled Mate-Pair libraries were sequenced on an Illumina HiSeq 2500 in Rapid mode (v1) with 2 × 101 bp paired-end reads. The loading concentration was 9.5 pM. This Initial run was for library QC purposes prior to running high throughput.

The Mate-Pair libraries were also sequenced on an Illumina HiSeq 2500 in high output mode (v4) with 2 × 125 bp paired-end reads. Libraries were sequenced on individual lanes (not pooled). The template loading concentration for each lane was adjusted based on the cluster density from the QC run. Two replicate flowcells were sequenced simultaneously, each with 6 lanes of mate-pair libraries.

Illumina read clouds (synthetic long reads) WGS

Library preparation. Synthetic long-read libraries were generated for the AJ Trio and Chinese Trio using the TruSeq Synthetic Long-Read DNA Library Prep Kit (Illumina, Cat# FC-126–1001). 500 ng of DNA from the NIST Reference Materials (or from Coriell for the Chinese parents) was sheared, end-repaired, A-tailed, and adapters ligated before size-selecting 9–11 kb fragments according to the manufacturer's protocol (Illumina Part # 15047264 Rev. B). Each resulting library was then diluted and aliquoted across a 384-well plate to limit the number of molecules to be amplified by PCR in each well. Amplified products were then tagged and indexed by a second round of PCR (see referenced protocol for conditions) before pooling and concentrating all 384 wells for final product size selection and validation, again according to manufacturer's instructions.

Sequencing. The synthetic long-read libraries for each genome were pooled and sequenced on an Illumina HiSeq 2500 in high output mode (v4) with 2 × 125 bp paired-end reads. Pooled libraries from each genome were loaded on individual lanes, with two lanes of each genome sequenced. The loading concentration for each lane was adjusted based on the cluster density of a previous (failed) run.

Illumina paired end WES

Library preparation. Captured exome DNA libraries for the Ashkenazim Jewish trio and the Chinese son were prepared from 4 vials of NIST Reference Material DNA. The library preparation was performed using Agilent SureSelect Automated Library Prep and Capture System Protocol and the Agilent SureSelect Target Enrichment System kit for 96 reactions (Agilent Technologies).

The DNA concentration for each sample was measured by Qubit 2.0 fluorometer (Life Technologies). The DNA samples were normalised with Low TE buffer to 4 ng/ul. 50 ul of each sample was fragmented in a 130 ul 96 MicroTUBE Plate with 96 microTUBE Foil Seal (part numbers 520078 and 520073 respectively) using the Covaris E-Series E220 focused ultrasonicator. The samples were sheared to an average of 250 bp, using the following instrument shearing settings: Average fragment size: 250 bp; Acousty Duty Factor: 10%; PeakIncidence Power, W: 140; Cycles Per Burst: 200; Treatment Time: 160 s; Temperature: 6 °C.

After fragmentation, the sheared DNA was purified using Agencourt AMPure XP purification beads, and the fragment size was controlled using the TapeStation 2200 with D1000 Reagents and ScreenTape. Following purification, the sample ends were modified for SureSelect target enrichment, through end-repair, 3' end adenylation, and adaptor ligation. The prepared DNA was purified after each each modification step. Sample purification, GA end-repair, A-tailing, and adaptor ligation were performed on the Bravo Agilent NGS Workstation robot in Nunc DeepWell plates. The adapter-ligated DNA fragments were captured and amplified in a 0.2 ml ABGENE 96-well PCR plate using the Applied Biosystems Veriti 96-well thermocycler. The thermocycler was programmed with the following settings: 98 °C/3 min, 10 cycles of (98 °C/80 s, 65 °C/30 s, 72 °C/1 min), 72 °C/10 min, and 4 °C/hold.

The libraries were cleaned up in a Nunc DeepWell plate using Agencourt AMPure XP beads on the Agilent NGS Workstation. To verify the DNA fragments has a size distribution between 200 and 400 bp, the libraries were measured on the TapeStation 2200 with D1000 Reagents and ScreenTape. The library concentration was measured using the Quant-iT dsDNA High-Sensitivity Assay Kit with the Victor3 1420 Multilabel counter. 750 ng of each prepped library sample was aliquoted for hybridization to the SureSelect Capture Library. Hybridization of the DNA libraries and the SureSelect Capture libraries was performed in a Nunc DeepWell plate on the Agilent NGS Workstation, and in a 0.2 ml ABGENE 96-well PCR plate on the Applied Biosystems Veriti 96-well thermocycler. The thermocycler was programmed with a 95 °C hybridization step for 5 min followed by a 65 °C hold step. The hybridized libraries were captured with SureSelect Binding Buffer and purified using Dynabeads MyOne Streptavidin T1 bead suspension.

Addition of indexing tags to the SureSelect enriched captured libraries was performed through PCR based amplification, using the Agilent NGS Workstation and Applied Biosystems Veriti 96-well thermocycler. The thermocycler was programmed with the following settings: 98 °C/2 min, 10 cycles of (98 °C/30 s, 57 °C/30 s, 72 °C/1 min), 72 °C/10 min, and 4 °C/hold.

The amplified and indexed DNA libraries were purified with Agencourt AMPure XP beads. Quality control of the amplified captured libraries was done using the TapeStation 2200 High Sensitivity D1000 Reagents and ScreenTape to ensure a fragment size of 300–400 bp and a concentration of 10 nM. To estimate the cluster density and achieve the final library concentration, qPCR was performed for each library using the KAPA library quantification kit for Illumina sequencing platforms. All four samples were prepared equally.

Sequencing. The libraries were pooled and diluted to a 2 nM pool (0.5 nM from each library), based on the qPCR measures. The library pool was sequenced on the Illumina HiSeq 2500 sequencing platform in high output run mode with 2×125 bp paired-end reads.

Analysis. The sequencing data were aligned by `bwa mem`⁶ against b37 human decoy reference genome. The alignments were sorted and PCR duplicates were marked by Picard (<http://picard.sourceforge.net>). For AJ trios, a joint variant calling was performed by GATK⁷ HaplotypeCaller on all three samples. For the Chinese son, both single sample variant calling (a VCF file) and the first step in cohort analysis (a gVCF file) were performed by GATK HaplotypeCaller. All variants in VCF files were quality filtered by standard GATK SNP variant quality score recalibration and indel hard filtration according to GATK Best Practices recommendations^{8,9}.

10X Genomics GemCode libraries for illumina sequencing

Genomic DNA extraction. Genomic DNA was purified using a modified version of the MagAttract HMW DNA Kit (QIAGEN, Germantown, MD) from GM12878, GM24149, GM24143 and GM24385 cells (Coriell, Camden, New Jersey). Briefly, 1×10^6 cells per extraction were pelleted and washed with PBS at RT. The Proteinase K and RNaseA digestion was incubated for 30 min at 25 °C. Genomic DNA was purified using MagAttract Suspension G with Buffer MB, washed twice with Buffer MW1, and twice with Buffer PE. Finally the beads were rinsed twice with nuclease-free water for exactly 60 s. DNA was eluted with Buffer AE and quantified using the Qubit dsDNA HS Assay Kit (Thermo Fisher Scientific, Waltham, MA).

GemCode whole genome library preparation and sequencing. Sample indexed and partition barcoded libraries were prepared using the GemCode kit (10X Genomics, Pleasanton, CA). 1.2 ng of DNA was used for GEM reactions where DNA fragments were massively partitioned into molecular reactors to extend the DNA and introduce specific 14-bp partition barcodes. GEM reactions were thermal cycled (95 °C for 5 min; cycled 18X: 4 °C for 30 s, 45 °C for 1 s, 70 °C for 20 s, and 98 °C for 30 s; held at 4 °C) and purified using the GemCode protocol. Purified DNA was sheared to 800-bp (M220, Covaris, Woburn, MA). Peak incident power: 75.0 W; duty factor: 5.0%; cycles per burst: 200; treatment time: 50 (s), temperature: 20.0 °C; sample volume 50 ul. End repair, Adenylation tailing of 3' ends, universal adapter ligation and sample indexing were performed according to the manufacturer's recommendations. Whole genome GemCode libraries were quantified by qPCR (KAPA Library Quantification Kit for Illumina platforms, Kapa Biosystems, Wilmington, MA). The NA12878 library was pooled with the NA24149 library and run on an Illumina HiSeq 2500 in Rapid mode (v1) with paired end 2×98 -bp, 14-bp I5 and 8-bp I7 reads. For analysis the demultiplexed results from three flow cells were combined for a total of approximately 1.25 billion and approximately 810 million reads for NA12878 and NA24149, respectively. The NA24385 and NA24143 libraries were each run individually on a single high output mode (v4) lane for approximately 958 and approximately 900 million reads, respectively. Sequencing results were analyzed using the GemCode Long Ranger Software Suite.

Complete Genomics WGS

Library preparation. Except for the Chinese parents, the NIST reference material was used as input. In the case of the Chinese trio, the full trio was sequenced from cells purchased from Coriell, and the son (GM24631) was also sequenced from the NIST reference material. Library prep followed the basic approach detailed previously¹⁰, but with a two adapter library protocol (library version 2). Briefly, sequencing substrates were generated by means of genomic DNA fragmentation to a median fragment length of about 450 base pairs and recursive directional adapter insertion with an intermediate type IIS restriction enzyme digestion. The resulting circles were then replicated with *Phi*29 polymerase (RCR)¹¹ by synchronized synthesis to obtain hundreds of tandem copies of the sequencing substrate, referred to as DNA nanoballs (DNBs) which were adsorbed to silicon substrates with grid-patterned arrays to produce DNA nanoarrays.

Sequencing. High-accuracy cPAL sequencing chemistry (Version 2 sequencing) was used on automated sequencing machines to independently read up to 19 bases adjacent to each of the four anchor insertion sites, resulting in a total of 29-base mate-paired reads (58 bases per DNB). DNB intensity information is interpreted with the following steps: 1) background removal, 2) image registration, 3) intensity extraction. The intensity data from each field were then subjected to base calling, which involved four major steps: 1) crosstalk correction, 2) normalization, 3) base calling, and 4) raw base score computation.

Complete Genomics LFR

Library preparation. The LFR libraries were constructed as described in ref. 12, except using the two adapter library protocol (library version 2) described above. Because this protocol requires cells as input, Coriell cells were used rather than the NIST reference material DNA. Briefly, controlled random enzymatic fragmenting is applied to 100–130 pg of high molecular mass (HMM) DNA that is physically separated into 384 distinct wells. The resulting fragments are then amplified and ligated to uniquely barcoded adapters. After combining the 384 wells and performing a restriction digestion, the second adapter is attached. The resulting substrate is converted to DNBs and adsorbed to silicon substrates with grid-patterned arrays to produce DNA nanoarrays.

Sequencing. Sequencing was performed as described above for regular Complete Genomics WGS, with the additional step of sequencing the well ID barcodes.

Ion exome sequencing

Library preparation. Exome libraries for 4 NIST Reference Materials, the AJ trio and Chinese son, were prepared using Ion AmpliSeq Exome RDY Kit, with a mean insert size of 215 bp. Each sample was assigned a distinct barcode: IonXpress_020 for NA24385, IonXpress_022 for NA24149, IonXpress_024 for NA24143, and IonXpress_026 for NA24631. Each barcode library is diluted to 100 pM. The libraries were emulsion-amplified individually and enriched using Ion OneTouch 2 System and Ion PI Template OT2 200 Kit v4. Outputs from 4 OneTouch runs for each sample were pooled together.

Sequencing. Each sample was sequenced on 4 Ion Proton instruments using Ion PI Sequencing 200 Kit v4. BaseCalling and alignment were performed on a Torrent Suite v4.2 server.

SOLiD 5500×1 Wildfire WGS

Fragment library preparation and sequencing of AJ son and Chinese son. DNA sequencing on a Life Technologies 5500×1 Wildfire was performed according to manufacturers protocols with noted modifications for each genome. A semi-automated library preparation process was first performed for the Chinese son Reference Material DNA. A modified manual library preparation was performed for the AJ son Reference Material DNA in an attempt to obtain smaller libraries for AJ son to maximize efficiency of colony formation on the 5500 W. The two procedures used for each genome are detailed below.

Chinese son semi-automated library preparation. A semi-automated library preparation using the AB Library Builder System was used to prepare Chinese son libraries for sequencing on a 5500×1 Wildfire. The workflow to produce 5500 W DNA fragment libraries from Chinese son human genomic DNA (gDNA) was as follows (also in Supplementary Fig. 1):

Shearing of gDNA was performed using the Covaris g-Tube (PN 520079) in conjunction with the Covaris S2 Focused Ultrasonicator. To obtain a uniform intermediate size distribution of approximately 10 kb, 2.0 ug of gDNA was initially ‘pre-sheared’ using a Covaris g-Tube in an Eppendorf 5424 centrifuge. The g-tubes were centrifuged twice at 4200 rpm for 60 s in each direction. Shearing was completed using the Covaris S2 per the User Guide ‘Fragment Library Preparation Using the AB Library Builder System: 5500 Series SOLiD Systems’ (PN 4460965 Rev. A). Approximately 1.5 ug of ‘pre-sheared’ gDNA was further sheared using the Covaris S2. Shearing was assessed on an Agilent 2100 Bioanalyzer High Sensitivity DNA Chip (PN 5067–4626) which showed a broad distribution of sheared material with peak at approximately 175 bp.

The Life Technologies AB Library Builder System was used to partially automate the library preparation process. End Repair, Size Selection, PolyA Tailing and Adaptor Ligation were performed on the AB Library Builder System to generate 5500 DNA fragment libraries. The Life Technologies Library Builder Fragment Core Kit for 5500 Genetic Analysis Systems (PN 4463763) and Beckman Coulter Agencourt AMPure XP Reagent (PN A263800) were used to prepare 5500 libraries on the AB Library Builder System. Adaptor amounts were calculated, per the Library Builder User Guide, based on input mass for a given sample. Library Preparation input mass ranged from 1.0–1.5 ug of sheared DNA depending on the given sample.

The AB Library Builder 5500 libraries then underwent manual nick translation and Wildfire library conversion to prepare libraries compatible for sequencing on a Life Technologies 5500×1 Wildfire. Wildfire conversion was performed per the Quick Reference ‘5500 W Series Genetic Analysis Systems: Conversion of 5500 Library to 5500 W Library’ (PN 4477188 Rev. B). Six cycles of amplification were performed in the conversion process.

Following an AmPure XP Reagent Cleanup the final 5500 W DNA fragment libraries were run on the Sage Science BluePippin automated DNA size selection and collection system to further narrow the size distribution of the final libraries. A BluePippin DNA 2% Dye-Free Agarose gel cassette with V1 Marker (PN BDF2010) was used to capture DNA in a target range of 200–300 bps. All 5500 W library for a given sample was loaded into the assigned well on cassette and run per the BluePippin 2% Agarose Gel Cassette Quick Guide. Upon completion of size selection 40–60 ul of size selected library was removed from the elution well and cleaned and concentrated using a 1.8X Agencourt AMPure XP (PN A263800) cleanup. Cleaned-size selected DNA was eluted in 32 ul of TE buffer. Size selection assessed using a Bioanalyzer

High Sensitivity DNA Chip and showed the final Chinese son 5500 W libraries with a size distribution of approximately 200–350 bps with peak at approximately 285 bps.

AJ son modified manual library preparation. A modified manual library preparation process for the AJ son was used to obtain appropriately sized libraries for sequencing on a 5500×1 Wildfire. The workflow to produce 5500 W DNA fragment libraries from AJ son human genomic DNA (gDNA) was as follows (also in Supplementary Fig. 2):

Shearing of gDNA was performed using the Covaris g-Tube (PN 520079) in conjunction with the Covaris S2 Focused Ultrasonicator. To obtain a uniform intermediate size distribution of approximately 10 kbp, 2.5 ug of gDNA was initially sheared using a Covaris g-Tube in an Eppendorf 5424 centrifuge. The g-tubes were centrifuged twice at 4200 rpm for 60 s in each direction. Shearing was completed using the Covaris S2 per the User Guide for 'Fragment Library Preparation: 5500 Series SOLiD Systems' (PN 4460960 Rev. B). Approximately 2.0 ug of 'pre-sheared' gDNA was sheared using the Covaris S2.

DNA fragment library preparation was performed using the total mass of sheared DNA (approximately 2.0 ug). Following the aforementioned 5500 Fragment Library Preparation guide, the 5500 SOLiD Fragment Library Core Kit (PN 4464412) was used to prepare 5500 libraries. The ends of the DNA fragments were repaired and DNA was cleaned and concentrated. Prior to size selection, the fragmented end-repaired DNA was assessed on an Agilent 2100 Bioanalyzer High Sensitivity DNA Chip (PN 5067–4626). Shearing resulted in a broad distribution with peak at approximately 175 bps.

To obtain a narrow fragment size distribution of DNA for AJ son 5500 library preparation, the DNA was run on the Sage Science BluePippin automated DNA size selection and collection system. A BluePippin DNA 3% Dye-Free Agarose gel cassette with Q2 Marker (PN BDF310) was used to capture DNA in a target range of 100–150 bps. Approximately 1.0–1.5 ug of end-repaired DNA was loaded into an appropriate well on the cassette and run per the BluePippin 3% Agarose Gel Cassette Quick Guide. Upon completion of size selection 40–60 ul of size selected DNA was removed from the elution well and cleaned and concentrated using a 1.8X Agencourt AMPure XP (PN A263800) cleanup. Cleaned size-selected DNA was eluted in 32 ul of TE buffer. Size selection was again assessed using a Bioanalyzer High Sensitivity DNA Chip and showed a peak at approximately 128 bps.

Using reagents provided in the 5500 SOLiD Fragment Library Core Kit, dA Tailing, adaptor ligation and nick translation were performed per the 5500 Fragment Library Preparation Guide. Adaptor volumes were calculated using the mass calculated from the Bioanalyzer High Sensitivity Chip following the cleanup of the size-selected DNA. Two rounds of cleanup using Agencourt AMPure XP reagent were performed per the User Guide.

The completed 5500 libraries were then converted to 5500 W libraries compatible for sequencing on a Life Technologies 5500×1 Wildfire. Wildfire conversion was performed per the Quick Reference '5500 W Series Genetic Analysis Systems: Conversion of 5500 Library to 5500 W Library' (PN 4477188 Rev. B) utilizing reagents provided in the Life Technologies 5500 W Conversion Primer Kit (PN 4478020) and Platinum PCR SuperMix (PN 11306–081). Six cycles of amplification were performed in the conversion process. Following completion of the conversion process, cleaned and concentrated libraries were assessed using a Bioanalyzer High Sensitivity DNA Chip and showed a peak at approximately 270 bps with a distribution from approximately 150–400 bps.

A second round of BluePippin size selection was performed to tighten the size distribution of the final 5500 W library. The 5500 W libraries were run on a DNA 2% Dye-Free Agarose gel cassette with V1 Marker to capture DNA in a target range of 200–300 bps. All DNA for a given sample was loaded into the assigned well on the cassette and run per the BluePippin 2% Agarose Gel Cassette Quick Guide. Upon completion of size selection 40–60 ul of size selected DNA was removed from the elution well and cleaned and concentrated using a 1.8X Agencourt AMPure XP cleanup. Cleaned-size selected 5500 W libraries were eluted in 32 ul of TE buffer. Size selection was assessed using a Bioanalyzer High Sensitivity DNA Chip and showed the final AJ son 5500 W libraries with a peak at approximately 275 bps and a distribution from approximately 240–320 bps.

Sequencing of AJ son and Chinese son. A Life Technologies 5500×1 Wildfire (5500 W Genetic Analysis System) was used to sequence 5500 W AJ son and Chinese son libraries using ICS software version 2.1. The User Guide '5500 W Series Genetic Analysis System (Americas)' (PN 4481746 Rev. B) was followed and used to prepare the samples and load a 5500 W v2 FlowChip (PN 4475661). The Wildfire Template Amplification Protocol v6.1, located in the User Guide, was followed for template amplification. The 5500 W FlowChip Prep Enzyme Kit (PN 4481058) and 5500 W Template Amplification Reagents v2 (PN 4475663) were used to prepare FlowChips for 'on-instrument' template amplification following template hybridization. 5500 W library molar concentrations were calculated from the Bioanalyzer High Sensitivity chip following the final size selection of the 5500 W libraries. These concentrations were used in calculation of FlowChip loading concentrations. Libraries were deposited into individual lanes at final concentrations of 100 to 250 pM. The library concentrations vary due to adjustments in subsequent instrument runs to increase colony density for a given library on the FlowChip.

5500 W fragment libraries were sequenced with single-end 75 bp reads using the 5500 W Forward SR 75 Reagent (PN 4475685). Two libraries were prepared and sequenced for each genome for a total of 24 lanes (4 FlowChips) per genome. This sequencing yielded approximately 72x coverage/genome.

Bionano genomics genome maps

Library preparation. Lymphoid-cell lines from the AJ Trio cell cultures obtained from Coriell Cell Repositories (GM24385, GM24143 and GM24149) were pelleted and washed with Life Technologies PBS (phosphate-buffered saline) at 1X concentration; the final cell pellet was re-suspended in cell suspension buffer using the Bio-Rad CHEF Mammalian Genomic DNA Plug Kit. Cells were then embedded in Bio-Rad CleanCut low melt Agarose and spread into a thin layer on a custom support in development. Cells were lysed using BioNano Genomics IrysPrep Lysis Buffer, protease treated with QIAGEN Puregene Proteinase K, followed by brief washing in Tris with 50 mM EDTA and then washing in Tris with 1 mM EDTA before RNase treatment with Qiagen Puregene RNase. DNA was then equilibrated in Tris with 50 mM EDTA and incubated overnight at 4 °C before extensive washing in Tris with 0.1 mM EDTA followed by equilibration in New England BioLabs NEBuffer 3 at 1X concentration. Purified DNA in the thin layer agarose was labeled following the BioNano Genomics IrysPrep Reagent Kit protocol with adaptations for labeling in agarose. Briefly, 1.25 µg of DNA was digested with 0.7 units of New England BioLabs Nt.BspQI nicking endonuclease per µl of reaction volume in New England BioLabs NEBuffer 3 for 130 min at 37 °C, then washed with Affymetrix TE Low EDTA Buffer, pH 8.0, followed by equilibration with New England BioLabs 1x ThermoPol Reaction Buffer. Nick-digested DNA was then incubated for 70 min at 50 °C using BioNano Genomics IrysPrep Labeling mix and New England BioLabs Taq DNA Polymerase at a final concentration of 0.4 U/µl. Nick-labeled DNA was then incubated for 40 min at 37 °C using BioNano Genomics IrysPrep Repair mix and New England BioLabs Taq DNA Ligase at a final concentration of 1 U/µl. Labeled-repaired DNA was then recovered from the thin layer agarose by digesting with GELase and counterstained with BioNano Genomics IrysPrep DNA Stain prior to data collection on the Irys system.

DNA was isolated from a lymphoid-cell culture of the Chinese son (GM24631) using the Bio-Rad CHEF Mammalian Genomic DNA Plug Kit protocol and lysed using BioNano Genomics IrysPrep Lysis Buffer and digested with QIAGEN Puregene Proteinase K. DNA was solubilized using GELase Agarose Gel-Digesting Preparation and drop-dialyzed before labeling using standard IrysPrep Reagent Kit protocols.

Mapping. Labeled and stained DNA samples were loaded into BioNano Genomics IrysChips and run on the BioNano Genomics Irys System imaging instrument. Data was collected for each sample until desired fold coverage of long molecules (>150 kb) was achieved. BioNano Genomics IrysView visualization and analysis software application was used to detect individual linearized DNA molecules using the Life Technologies YOYO-1 Iodide in DMSO and determine the localization of labeled nick sites along each DNA molecule. BioNano Genomics IrysSolve analytical and assembly pipeline compiled the sets of single-molecule maps for each sample and were then used to build a full genome assembly.

Pacific biosciences

SMRTbell library preparation of AJ Trio gDNA. DNA library preparation and sequencing was performed according to the manufacturer's instructions with noted modifications. Following the Pacific Biosciences Protocol, '20-kb Template Preparation Using Blue Pippin Size-Selection System', library preparation was performed using the Pacific Biosciences SMRTbell Template Prep Kit 1.0 (PN # 100-259-100). In short, 10 µg of extracted, high-quality, genomic DNA from the NIST Reference Material DNA for the AJ trio, were used for library preparation. Genomic DNA extracts were verified with the Life Technologies Qubit 2.0 Fluorometer using the High Sensitivity dsDNA assay (PN# Q32851) to quantify the mass of double-stranded DNA present. After quantification, each sample was diluted to 150 µl, using kit provided EB, yielding a concentration of approximately 66 ng/µl. The 150 µl aliquots were individually pipetted into the top chambers of Covaris G-tube (PN# 520079) spin columns and sheared for 60 s at 4500 rpm using an Eppendorf 5424 benchtop centrifuge. Once complete, the spin columns were flipped after verifying that all DNA was now in the lower chamber. The columns were spun for another 60 s at 4500 rpm to further shear the DNA and place the aliquot back into the upper chamber. In some cases G-tubes were centrifuged 2–3 times, in both directions to ensure all volume had passed into the appropriate chamber. Shearing resulted in a approximately 20,000 bp DNA fragments verified using an Agilent Bioanalyzer DNA 12000 gel chip (PN# 5067–1508). The sheared DNA isolates were then purified using a 0.5X AMPure PB magnetic bead purification step (0.5X AMPure PB beads added, by volume, to each DNA sample, vortexed for 10 min at 2,000 rpm, followed by two washes with 70% alcohol and finally eluted in EB). This AMPure purification step assures removal of any small fragment and/or biological contaminant. The sheared DNA concentration was then measured using the Qubit High Sensitivity dsDNA assay. These values were used to calculate actual input mass for library preparation following shearing and purification.

After purification, approximately 8 to 9 µg of each purified sheared sample went through the following library preparation process per this protocol (also in Supplementary Fig. 3):

All library preparation reaction volumes were scaled to accommodate input mass for a given sample. Library size selection was performed using the Sage Science BluePippin 0.75% Agarose, Dye Free, PacBio approximately 20 kb templates, S1 cassette (PN# PAC20KB). Size selections were run overnight to maximize recovered mass. Approximately 2–5 mg of prepared libraries were size selected using a 10 kb start and 50 kb end in ‘Range’ mode. This selection is necessary to narrow the library distribution and maximize the SMRTbell sub-read length for the best *de novo* assembly possible. Without selection, smaller 2000–10,000 bp molecules dominate the zero-mode waveguide loading distribution, decreasing the sub-read length. Size-selection was confirmed using pre and post size selected DNA using an Agilent DNA 12000 chip. Final library mass was measured using the Qubit High Sensitivity dsDNA Assay. Approximately 15–20% of the initial gDNA input mass resulted after elution from the agarose cassette, which was enough yield to proceed to primer annealing and DNA sequencing on the PacBio RSII instrument. This entire library preparation and selection strategy was conducted 7, 2 and 2 times across AJ son, AJ father, and AJ mother respectively, to provide enough library for the duration of this project.

Sequencing AJ trio on pacific Biosciences RSII. Sequencing reflects the P6-C4 sequencing enzyme and chemistry, respectively. (Note that 10.3% of the data was collected using the P5-C3 enzyme/chemistry prior to the release of the P6-C4 enzyme and chemistry.) Primer was annealed to the size-selected SMRTbell with the full-length libraries (80 °C for 2 min 30 followed by decreasing the temperature by 0.1°/s to 25°C). To prepare the polymerase-template complex, the SMRTbell template complex was then bound to the P6 enzyme using the Pacific Biosciences DNA Polymerase Binding Kit P6 v2 (PN# 100-372-700). A ratio of 10:1, polymerase to SMRTbell at 0.5 nM, was prepared and incubated for 4 h at 30 °C and then held at 4 °C until ready for magbead loading prior to sequencing. The Magnetic bead-loading step was conducted using the Pacific Biosciences MagBead Kit (PN# 100-133-600) at 4 °C for 60-minutes per manufacturer’s guidelines. The magbead-loaded, polymerase-bound, SMRTbell libraries were placed onto the RSII instrument at a sequencing concentration of 100 to 40 pM to optimize loading across various SMRTcells. Sequencing was performed using the C4 chemistry provided in the Pacific Biosciences DNA Sequence Bundle 4.0 (PN# 100-356-400). The RSII was then configured for at least 240-minute continuous sequencing runs.

Oxford Nanopore

The genomic DNA library preparation consists of the ligation of a hairpin adapter to dsDNA molecules (either sheared to ~8 kb as is currently recommended for optimal data yield or unsheared to produce the longest possible reads) such that the template, then adapter, then complement can be passed through the pore sequentially. This library design produces a current time-series dataset with three distinct sections, of which the template and complement can be isolated from the adapter region. After base-calling is performed, the template and complement are aligned to produce two-direction, or 2D, reads. If the quality of one or both sequences is limited, a 2D read may not be produced, though a 1D read is made available.

Library preparation of AJ Son gDNA. Genomic DNA from the Ashkenazi Jewish (AJ) son was prepared for sequencing via the Oxford Nanopore Technologies MinION single molecule sequencing instrument. Two libraries were generated, one with the ‘SQK-MAP-004 genomic DNA’ kit and one with the ‘SQK-MAP-006 genomic DNA’ kit provided as part of the MAP. Library preparation and sequencing was done according to manufacturer’s instructions with all optional steps executed. Both libraries were prepared with 1 µg HMW-gDNA of the HG-002 RM. DNA concentration was measured using Life Technologies Qubit dsDNA BR assay (PN# Q32850). DNA quality was measured with the Agilent 2200 TapeStation Genomic DNA Analysis assay (PN# 5067–5365). Shearing was done with Covaris G-tubes (PN#520079) and an Eppendorf 5424R centrifuge (PN# 5404000413). Prior to library preparation, the optional New England BioLabs preCR repair (PN# M0309S) step was taken for the SQK-MAP-004 library and the corresponding optional NEB FFPE (PN# M6630S) repair step was taken for the SQK-MAP-006 library. The RM was never stored via FFPE but the protocol distributed by ONT suggests using this particular reagent prior to library preparation to repair potential DNA damage in the interest of producing the highest quality signal during sequencing.

Sequencing of AJ Son on Oxford Nanopore MinION. Both AJ Son gDNA libraries were sequenced via single 48 h runs on the MinION instrument. A MinION version 7.3 flow cell was used for the SQK-MAP-004 library and a MinION MkI for the SQK-MAP-006 library, each being the most current version at the time of sequencing. Flowcells were received along with library preparation kits as part of the MAP. Flowcells were primed twice at 10 min intervals prior to loading the library, as described in each respective library preparation protocol. Sequencing runs were controlled using default versions of MinKNOW protocols ‘MAP_48Hr_Sequencing_Run.py’ (SQK-MAP-004 library) and ‘MAP_48Hr_Sequencing_Run_SQK_MAP006.py’ (SQK-MAP-006 library).

Methods to generate 2D reads from the Oxford Nanopore MinION from two different libraries are described in ref. 13.

Data Records

Genomic samples

The genomes sequenced in this work (see Table 2) and their data are all publicly available both as EBV-immortalized B lymphoblastoid cell lines (from Coriell only) and as DNA (from Coriell and NIST). As described in the Methods, most data are from the NIST Reference Materials, unless the technology benefited from preparing longer DNA directly from cells. These human subjects are approved for ‘Public posting of personally identifying genetic information (PIGI)’ by the Coriell and NIH/NIGMS IRBs, and this study was approved by NIST and by the Coriell/NIGMS IRB.

Illumina paired end WGS

NA12878. Approximately 300x 148 bp × 148 bp Illumina paired end WGS data from NA12878 (HG001) is in the NCBI SRA SRX1049768 to SRX1049855 [Data Citation 1].

AJ Trio. 148 × 148 bp HiSeq sequencing and analyses of two 40x to 50x runs from each member of the Ashkenazim trio (HG002, HG003, and HG004) using the BWA-GATK pipeline on Basespace. Raw data is available in the SRA: SRX847862 to SRX848317 [Data Citation 2]. The BAM, VCF, and fastq files have been uploaded to:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_HiSeq_HG002_Homogeneity-10953946/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_HiSeq_HG003_Homogeneity-12389378/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_HiSeq_HG004_Homogeneity-14572558/

Chinese Trio. Data from the Chinese trio is in the NCBI SRA SRX1388368 to SRX1388459 [Data Citation 3]. Fastq files for 300x sequencing of Chinese son (HG005), as well as approximately 45x bam files generated from each flow cells are located here:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/HG005_NA24631_son_HiSeq_300x

Fastq files for 100x sequencing of the Chinese parents (GM24694 and GM24695), as well as approximately 100x bam files generated for each genome are located here:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG006_NA24694-huCA017E_father/NA24694_Father_HiSeq100x

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG007_NA24695-hu38168_mother/NA24695_Mother_HiSeq100x

Illumina mate-pair sequencing

Illumina mate-pair data are available at the NCBI SRA SRX1388732 to SRX1388743 [Data Citation 4], and as bam files in the NIST_Stanford_Illumina_6 kb_matepair directory for each genome (HG002, HG003, HG004, HG005, GM24694, and GM24695):

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_Stanford_Illumina_6kb_matepair/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_Stanford_Illumina_6kb_matepair/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_Stanford_Illumina_6kb_matepair/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/NIST_Stanford_Illumina_6kb_matepair/

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG006_NA24694-huCA017E_father/NIST_Stanford_Illumina_6kb_matepair/

Genome	Coriell cell line ID	NIST ID	NIST RM #	NCBI BioSample	PGP ID
CEPH Mother/Daughter	GM12878	HG001	RM8398	SAMN03492678	Not PGP
AJ Son	GM24385	HG002	RM8391* (son)/RM8392* (trio)	SAMN03283347	huAA53E0
AJ Father	GM24149	HG003	RM8392* (trio)	SAMN03283345	hu6E4515
AJ Mother	GM24143	HG004	RM8392* (trio)	SAMN03283346	hu8E87A9
Chinese Son	GM24631	HG005	RM8393*	SAMN03283350	hu91BD69
Chinese Father	GM24694	N/A [†]	N/A [†]	SAMN03283348	huCA017E
Chinese Mother	GM24695	N/A [†]	N/A [†]	SAMN03283349	hu38168C

Table 2. Genome in a Bottle Consortium Genomes. Identifiers associated with genomes currently being characterized by GIAB. The NCBI BioProject for GIAB is PRJNA200694. *Not currently available. Planned release as NIST RMs in mid-2016. †NIST Reference Materials are not planned for the Chinese parents, but cells and DNA are available from Coriell.

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG007_NA24695-hu38168_mother/NIST_Stanford_Illumina_6kb_matepair/

Illumina read clouds (synthetic long reads)

Illumina read cloud data are available as fastq's (for the AJ trio and Chinese trio) and as bam files (currently only for the AJ son and father) in the NIST_Stanford_Moleculo directory for each genome (HG002, HG003, HG004, HG005, GM24694, and GM24695):

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_Stanford_Moleculo/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_Stanford_Moleculo/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_Stanford_Moleculo/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/NIST_Stanford_Moleculo/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG006_NA24694-huCA017E_father/NIST_Stanford_Moleculo/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG007_NA24695-hu38168_mother/NIST_Stanford_Moleculo/

Illumina paired end WES

Illumina paired-end WES data for AJ Trio (HG002, HG003, and HG004) and the Chinese son (HG005) are available. Raw data (fastq files) are available in the SRA: SRP047086 [Data Citation 5]. The BAM and VCF files have been uploaded to:

BAM files and their index files:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/OsloUniversityHospital_Exome
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/OsloUniversityHospital_Exome
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/OsloUniversityHospital_Exome
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/OsloUniversityHospital_Exome
 Joint Variant Calling file for AJ trio:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/OsloUniversityHospital_Exome_GATK_jointVC_11242015
 Single sample variant calling file for the Chinese son:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/analysis/OsloUniversityHospital_Exome_GATK_jointVC_11242015
 gVCF for the Chinese son:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/analysis/OsloUniversityHospital_Exome_GATK_jointVC_11242015

10X Genomics GemCode libraries for illumina sequencing

10X Genomics data was generated with cell lines acquired from Coriell (GM12878, GM24385, GM24149, and GM24143). Aligned reads with barcode and phasing information are provided in BAM format for each sample. VCF files with small variants are also provided for each sample, and SV calls are provided for NA12878 and the AJ son. See <http://software.10xgenomics.com/> for detailed information on file formats. 10X Genomics data are available at <http://software.10xgenomics.com/giab2015>

The same data are available at the NCBI SRA SRX1392293 to SRX1392296 [Data Citation 6] and on the GIAB FTP:

<ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/10XGenomics>
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/10XGenomics
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/10XGenomics
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/10XGenomics

Complete Genomics WGS

The Complete Genomics files are available on the NCBI SRA [Data Citation 7] and [Data Citation 8], and on the GIAB FTP site as summarized in Table 3. For the son of the Chinese trio, data is available for both the NIST reference material (HG005) and cells sourced from Coriell (GM24631). Data is available for the Chinese parents (GM24695 and GM24694) from cells sourced from Coriell. All other data is from NIST reference materials.

Directory structures and file formats for the ‘Full package’, as well as the other supplementary files discussed below, are described in http://www.completegenomics.com/documents/DataFileFormats_Standard_Pipeline_2.5.pdf.

For both the Chinese and Askenazi trios, a multisample VCF including only small variants was generated from masterVar files using the CGA tools mkvcf program described in <http://cgatools.sourceforge.net/docs/1.8.0/cgatools-user-guide.pdf>.

All other VCF files contain small variants, CNVs, SVs and MEIs. Note that for CNVs and SVs, more complete information is available in the ASM/CNV and ASM/SV directories of the full package. For CNVs, VCF files contain the information sourced from the cvnDetails files but do not provide information on any segmentation of the genome into ploidy or coverage levels. For SVs, VCF files contain information sourced from the allJunctionsBeta and highConfidenceJunctionsBeta, but information from the allSvEventsBeta and highConfidenceSvEventsBeta files is not included.

BAM files are provided in order to provide evidence of variants called. However, it is not appropriate to remap and recall variants based on these BAM files as proper re-mapping of reads should take into account the gapped read structure. The *_mapping_sorted_header.bam files include the initial mappings of all reads. They were generated with the map2sam program from CGA Tools with the --mate-sv-candidates and --add-unmapped-mate-info parameters. Inconsistent mappings are normally converted as single arm mappings with no mate information provided, but with the --mate-sv-candidates option map2sam will mate unique single arm mappings in SAM including those on different stands and chromosomes. The tag ‘XS:i:1’ is used to distinguish these ‘artificially’ mated records. The MAPQ provided for these records is a single arm mapping weight. The --add-unmapped-mate-info parameter generates mate sequences and score tags for inconsistent mappings. In the subsequent local *de novo* assembly (LDN) stage of genome assembly, reads can be re-mapped, added or removed from the assembly within the region undergoing LDN. The reads and the mappings that support variant calls after LDN is complete are provided in the evidence files. EvidenceDnbs* bam files are generated with our evidence2sam tool from CGA Tools (<http://cgatools.sourceforge.net/docs/1.8.0/cgatools-user-guide.pdf>).

Genome	File types	Location
HG001	VCF	/data/NA12878/analysis/variant_calls/COMPLETE/Library2/
HG001	BAM	/data/NA12878/CompleteGenomics_normal_RMDNA/alignment/Library2/
HG001	Full package	/technical/complete_data/HG001_NA12878_normalCG/GS000025639-ASM/EXP/GS02256-DNA_A01/
Ashkenazi trio	Multi-sample VCF	/data/AshkenazimTrio/analysis/CompleteGenomics_RefMaterial_SmallVariants_CGatools_08082014/AshkenaziTrio_RefMaterial_SmallVariants.vcf
HG002	BAM, VCF, reference scores	/data/AshkenazimTrio/HG002_NA24385_son/CompleteGenomics_normal_RMDNA/son_NA24385_GS000037263-ASM/
HG002	Full package	/technical/complete_data/trios/native_format_data/GS000043892-DID/GS000037263-ASM/
HG004	BAM, VCF, reference scores	/data/AshkenazimTrio/HG004_NA24143_mother/CompleteGenomics_normal_RMDNA/mom_NA24143_GS000037262-ASM/
HG004	Full package	/technical/complete_data/trios/native_format_data/GS000043891-DID/GS000037262-ASM/
HG003	BAM, VCF, reference scores	/data/AshkenazimTrio/HG003_NA24149_father/CompleteGenomics_normal_RMDNA/dad_NA24149_GS000037264-ASM/
HG003	Full package	/technical/complete_data/trios/native_format_data/GS000043893-DID/GS000037264-ASM/
Han trio (Coriell cells)	Multi-sample VCF	/data/ChineseTrio/analysis/CompleteGenomics_HanTrio_ExtractedFromCoriellCells_SmallVariants_CGatools_08082014/HanTrio_ExtractedFromCoriellCells_SmallVariants.vcf
HG005	BAM, VCF, reference scores	/data/ChineseTrio/HG005_NA24631_son/CompleteGenomics_normal_RMDNA/son_NA24631_GS000037265-ASM/
HG005	Full package	/technical/complete_data/trios/native_format_data/GS000043894-DID/GS000037265-ASM/
GM24631 (cells)	BAM, VCF, reference scores	/data/ChineseTrio/HG005_NA24631_son/CompleteGenomics_normal_cellsDNA/son_NA24631_GS000037475-ASM/
GM24631 (cells)	Full package	/technical/complete_data/trios/native_format_data/GS000044335-DID/GS000037475-ASM/
GM24695 (cells)	BAM, VCF, reference scores	/data/ChineseTrio/HG007_NA24695-hu38168_mother/CompleteGenomics_normal_cellsDNA/mom_NA24695_GS000037477-ASM/
GM24695 (cells)	Full package	/technical/complete_data/trios/native_format_data/GS000044337-DID/GS000037477-ASM/
GM24694 (cells)	BAM, VCF, reference scores	/data/ChineseTrio/HG006_NA24694-huCA017E_father/CompleteGenomics_normal_cellsDNA/dad_NA24694_GS000037476-ASM/
GM24694 (cells)	Full package	/technical/complete_data/trios/native_format_data/GS000044336-DID/GS000037476-ASM/

Table 3. Complete Genomics data locations. Locations of regular Complete Genomics data available at <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/>.

A detailed description of the data file can be found in the, ‘Representation of the Complete Genomics Data in SAM Output Format’ appendix of the CGA Tools User Guide (<http://cgatools.sourceforge.net/docs/1.8.0/cgatools-user-guide.pdf>). They contain the reads and mappings that support one of the called alleles by at least 2 dB over the other called allele. This means they will not contain reads and mappings that do not support either of the called alleles. The evidence BAMs do not contain reads and mappings for loci that were ultimately no-called or called homozygous ref, unless those regions were selected for *de novo* assembly because they were suspected to contain a variation. Every read that is found in the evidence files will also be present in the initial mappings, but the mapping positions may be different. In this case, where a read is found in both the *_mapping_sorted_header.bam and evidenceDnbs* files, the mapping in the evidence files is preferred.

Complete Genomics LFR

The Complete Genomics LFR data was sequenced from cells sourced from Coriell (GM12878, GM24385, GM24149, and GM24143). VCF and var formats files that include small variant calls are available for GM12878 and the Ashkenazi trio as indicated in Table 4. Both file formats include phasing information; see http://www.completegenomics.com/documents/DataFileFormats_Standard_Pipeline_2.5.pdf for details on file formats. In addition, summary files are included with assembly statistics.

The VCF and var files also include two additional FORMAT fields: MEWC and SWC. MEWC, minimum exclusive well count, indicates the number of LFR wells that support the REF or ALT allele (whichever is fewer) exclusively, and not the other allele. SWC, shared well count, indicates the number of LFR wells that support both the REF and ALT allele. High confidence variant calls should have a high MEWC (typically greater than 3) and a low (ideally 0) SWC. Note that small variant sensitivity is somewhat lower for the LFR process compared to standard Complete Genomics sequencing, so the standard sequencing should be deferred to for unphased variants.

Ion exome sequencing

The files generated by Thermo Fisher Scientific describe genomic variants called from Ion Torrent sequencing data with AmplicSeq exomes. The variants are represented in VCF files, each accompanied by an effective region BED file describing the corresponding genomic scope of called variants.

Four GIAB samples with were sequenced using AmpliSeq exome and sequenced on the Ion Proton (ThermoFisher Scientific, Catalog No. 4487084).

AmpliseqExome.20141120.16runs.vcf.zip -- 16 VCF files produced by Torrent Variant Caller v4.4 on 16 AmpliSeqExome runs, 4 of each samples picked by Genome in a Bottle consortium (HG002, HG003, HG004, and HG005).

AmpliseqExome.20141120.NA24143.vcf—HG004 variants called on 4 runs combined, above a quality score of 25;

AmpliseqExome.20141120.NA24149.vcf—HG003 variants called on 4 runs combined, above a quality score of 25;

AmpliseqExome.20141120.NA24385.vcf—HG002 variants called on 4 runs combined, above a quality score of 25;

AmpliseqExome.20141120.NA24631.vcf—HG005 variants called on 4 runs combined, above a quality score of 25;

AmpliseqExome.20141120_effective_regions.bed -- Genomic scope of AmpliSeqExome variant calls. This file describes the region in which variants are called with Torrent Suite v4.4 and later, and the region on which curation has been performed.

High_Confidence_Variants_NA24385.bed -- A list of inspected HG002 variants based on Ion Torrent, Complete Genomics, 23andme and manual curation

High_Confidence_Variants_NA24385_effective_regions.bed -- Genomic scope of inspected NA24385 variants.

The Ion Exome data are available on the NCBI SRA [Data Citation 9], [Data Citation 10], and [Data Citation 11], and on the GIAB FTP site at:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/ion_exome/
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/ion_exome/

GM12878 (rep 1)	VCF, var summary	/data/NA12878/CompleteGenomics_LFR/*GS000039392-ASM-NA12878*
GM12878 (rep 2)	VCF, var summary	/data/NA12878/CompleteGenomics_LFR/*GS000039396-ASM-NA12878*
GM12878 (rep 3)	VCF, var summary	/data/NA12878/CompleteGenomics_LFR/*GS000039473-ASM-NA12878*
GM24385	VCF, var summary	/data/AshkenazimTrio/HG002_NA24385_son/CompleteGenomics_LFR/*GS000039526-ASM-NA24385*
GM24143	VCF, var summary	/data/AshkenazimTrio/HG004_NA24143_mother/CompleteGenomics_LFR/*GS000039524-ASM-NA24143*
GM24149	VCF, var summary	/data/AshkenazimTrio/HG003_NA24149_father/CompleteGenomics_LFR/*GS000039541-ASM-NA24149*

Table 4. Complete Genomics LFR data locations. Locations of Complete Genomics LFR data available at <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/>.

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/ion_exome/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/ion_exome/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/ion_exome/

SOLiD 5500 × 1 Wildfire WGS

The SOLiD 5500 × 1 Wildfire WGS data for the AJ Son (HG002) and Chinese Son (HG005) are currently available as xsq files on the GIAB ftp site because this is the native format for SOLiD, as well as bam files mapped with lifescape and read groups for each lane. These data are available on the NCBI SRA [Data Citation 12] and [Data Citation 13], and on the GIAB FTP site at:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_SOLiD5500W
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/NIST_SOLiD5500W

Bionano Genomics genome maps

all.bnx is the raw data after image processing and filtering for molecules >150 kb

EXP_REFINEFINAL1.cmap is the *de novo* assembly consensus genome map set

The following files result from the alignment of genome maps to hg19:

EXP_REFINEFINAL1.xmap is the alignment file with match group information

EXP_REFINEFINAL1_q.cmap is the *de novo* genome maps that align to hg19 (query, it's a subset of all genome maps)

EXP_REFINEFINAL1_r.cmap is an in silico map of hg19 (Nt.BspQI motifs, anchor in the alignment)

BioNano data for the AJ Trio and Chinese Son (GM24385, GM24149, GM24143, and GM24631) are available at

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/BioNano/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/BioNano/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/BioNano/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/BioNano/

Pacific Biosciences

The PacBio data for the AJ Trio (HG002, HG003, and HG004) are available on the NCBI SRA [Data Citation 14] and on the GIAB FTP site at:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/PacBio_MtSinai_NIST/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG003_NA24149_father/PacBio_MtSinai_NIST/
 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG004_NA24143_mother/PacBio_MtSinai_NIST/

The file/directory naming convention is defined as follows: [SampleName]/[WellName]_[Collection-Number].[UUID].tar.gz Note that SampleName may contain other genomes in the name since this is hardcoded by the run name, but the data directories only contain run data from AJ son, AJ father, and AJ mother. For example, for SampleName of HG002new_O1_BP_P6_021815_MB_105 pM, WellName of A01, and CollectionNumber of 3, you will see a tar.gz file in HG002new_O1_BP_P6_021815_MB_105 pM directory with name A01_3.[UUID].tar.gz The UUID is currently used for only hashing purpose. The tar.gz file contains the raw SMRTPortal data including following contents:

```
tar.gz
| [movie name].1.xfer.xml
| [movie name].2.xfer.xml
| [movie name].3.xfer.xml
| [movie name].mcd.h5
| [movie name].metadata.xml
\---Analysis_Results
| [movie name].1.bax.h5
| [movie name].1.log
| [movie name].1.subreads.fasta
| [movie name].1.subreads.fastq
| [movie name].2.bax.h5
| [movie name].2.log
| [movie name].2.subreads.fasta
| [movie name].2.subreads.fastq
| [movie name].3.bax.h5
| [movie name].3.log
| [movie name].3.subreads.fasta
| [movie name].3.subreads.fastq
| [movie name].bas.h5
| [movie name].sts.csv
| [movie name].sts.xml
```

The metadata.xml contains all the metadata of this particular sample in the xml format; for example, in the TemplatePrep field you might see 'DNA Template Prep Kit 2.0 (3–10 Kb),' and in the BindingKit field you might see 'DNA/Polymerase Binding Kit P6,' etc. For information about bas.h5/bax.h5 files, please see: <http://files.pacb.com/software/instrument/2.0.0/bas.h5%20Reference%20Guide.pdf>

For information about subreads, please see: <https://speakerdeck.com/pacbio/track-1-de-novo-assembly>

Oxford Nanopore

The Oxford Nanopore raw reads and 2D reads for the AJ Son (HG002) are available at:

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/CORNELL_Oxford_Nanopore/

Future data and FTP structure

We expect to continue to accrue public data for these genomes as new methods become available. These data will be placed in the NCBI SRA when possible, linked to the GIAB BioProject PRJNA200694 and the appropriate BioSample listed in Table 2. Other data and analyses will also be publically available on the GIAB FTP site at NCBI (<ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/>). Preliminary data will be placed in the technical directory and analyses and finalized data will be placed under each trio or genome in the data directory. The directories under the analysis directory for each family contains the institution, data set, type(s) of variants, analysis tool, and date.

Technical Validation

Illumina paired end WGS

Several statistics were calculated for each flow cell using the Illumina BaseSpace Isaac Whole Genome Sequencing v3 analysis pipeline (Table 5).

Illumina mate-pair WGS

To assess duplication rate, coverage, and insert size of the mate-pair libraries, reads were stripped of adapter sequences. Read pairs were removed if the sequence of one or both mates was less than 20 bp after adapter stripping, or if the adapter sequence was at the beginning rather the end of a read (indicating the read inserts were likely to be in inward-facing F/R orientation rather than the expected outward-facing R/F orientation). Reads were then mapped to the hg19 reference genome using 'bwa mem'⁶ with default settings, and duplicates were marked using samblaster¹⁴. Statistics are summarized in Table 6.

The high rate of PCR duplicates (close to 50% in some libraries) resulted in lower than expected sequence coverage (13–17x average across all sequenced genomic positions). A more relevant metric for mate-pair data is the physical coverage, which measures the number of inferred fragments that cover a particular genomic position (including both the sequenced ends as well as the unsequenced genomic region between the ends). Because the empirical insert size average was between 6–7 kb per individual, the physical coverage of the genome was quite high (>400x per individual). BAMs were stripped of duplicate reads to reduce file size, but the full data are available in fastq format.

Illumina read clouds (synthetic long reads)

Several statistics were generated to assess the Illumina read clouds (Fig. 2): the read coverage distribution for each cloud, the fragment coverage distribution for the whole genome, the distribution for cloud length, and the probability density estimation for template length (same over samples). The cloud length distributions are plotted by type of the clouds. In such figures, '0 end' means both of the end-markers of a cloud are missing, and so on. Thus for clouds with both end-markers (2 ends), the expectation for the length is larger. Also for genome coverage distribution, the y-axis is the length of the genome covered by corresponding fragment coverage (integrate to 3 Gb).

Genome	Coverage (x)	Read length	percent duplicate paired reads	Fragment Length Median (bp)	Fragment length standard deviation (bp)	read 1 percent aligned	read 2 percent aligned	read 1 mismatch rate (%)	read 2 mismatch rate (%)
HG002	290	2 × 150	1.3	567	150	95.6	94.2	0.49	0.73
HG003	294	2 × 150	1.3	562	146	94.8	93.2	0.49	0.73
HG004	324	2 × 150	1.3	562	143	96.0	93.8	0.49	0.84
HG005	306	2 × 250	3.9	577	154	96.7	94.6	0.97	1.58
GM24694	105	2 × 150	3.7	560	146	95.8	93.9	0.51	0.78
GM24695	106	2 × 150	3.5	560	140	96.2	93.9	0.49	0.83

Table 5. Illumina paired end WGS statistics. Statistics were calculated using Illumina BaseSpace Isaac Whole Genome Sequencing v3 analysis pipeline and summarized by genome.

	after adapter stripping		uniquely mapping (mapq = 60, removing dups)	percent dups	sequence coverage		physical coverage		insert size	
	read count	base count (Gb)			mean	std	mean	std	mean	std
HG002	890,861,081	94.5	366,322,711	51	13.6	4.8	447	48.8	6466	1220
HG003	822,621,264	89.5	353,618,533	49	13.6	4.7	408	45.4	6110	934
HG004	815,116,039	88.5	379,255,921	45	14.2	4.9	421	49.4	6052	1375
HG005	663,052,149	71.4	382,392,628	32	14.3	5.1	463	51.7	6290	1982
GM24694	663,071,189	71.9	436,995,349	22	16.7	5.7	525	60.2	6517	1149
GM24695	648,828,312	71.1	418,948,777	24	15.8	5.5	507	64.9	6721	1204

Table 6. Illumina mate-pair sequencing statistics. PCR duplication rate was calculated to estimate library complexity, and insert size, sequence coverage (by reads), and physical coverage (by long fragments) were calculated after removing duplicates.

Illumina paired end WES

To assess duplication rate, coverage, and insert size of the libraries, picard CollectAlignmentSummaryMetrics, CollectInsertSizeMetrics and CalculateHsMetrics were performed on each sample BAM file. Statistics are summarized in Table 7.

10X Genomics GemCode libraries for illumina sequencing

Statistics for the AJ trio and Chinese son were generated using the 10X GemCode Long Ranger software (Table 8).

Complete Genomics WGS

Genomic assembly. Sequencing results in mate-paired reads with a 2–4 base overlap between adjacent contiguous sequences, as described in the ‘Read Data Format’ section of the Data File Formats documentation (http://www.completegenomics.com/documents/DataFileFormats_Standard_Pipeline_2.5.pdf). The gapped read pairs were aligned to the NCBI Build 37 reference genome using an index lookup based fast algorithm. At locations where the mapping results suggest the presence of a variant, mapped reads were refined, expanded and then assembled into a best-fit, diploid sequence with a custom software suite employing both Bayesian and de Bruijn graph techniques¹³. This process yielded diploid reference, variant or no-call at each genomic location with associated variant quality scores.

In addition to small variants, larger variants are detected, including MEIs, CNVs and SVs. The Complete Genomics CNV pipeline has the following steps: 1) Various measures of coverage for tiled 2 and 100 kb windows across the genome are determined. This is provided in the *cnvDetails** and *depthOfCoverage* files. 2) The genome is segmented into called ploidy levels (diploid model, 2 kb windows) or coverage levels (non diploid model, 100 kb windows) using diploid and non-diploid HMM-based algorithms. The segmentation patterns called are provided in the *cnvSegments** files. 3) The lesser allele fraction (LAF) is calculated for 100 kb windows across the genome—the LAF calculations are included in the *cnvDetailsNondiploid* and *cnvSegmentsNondiploid* files (and not in the diploid files because of the 100 kb window size restriction). To identify SVs, DNB mappings found during the standard assembly process are analyzed to find clusters of DNBs in which each arm maps uniquely to the reference genome, but with an unexpected mate pair length or anomalous orientation. SVs are encoded in the *junctions* and *highConfidenceJunctions* files where the latter file contains a high-confidence filtered subset of the data in former file. In addition to calling structural variant junctions, junctions are rationalized into structural variation events using the CGA Tools *junctions2events* algorithm. These data are provided in the *svEvents* and *highConfidenceSvEvents* files. Additional information on the CNV, SV and MEI algorithms is available here: http://www.completegenomics.com/documents/DataFileFormats_Standard_Pipeline_2.5.pdf

Assembly metrics are summarized in Table 9. Additional summary information can be found for each genome in the full package in the *ASM/summary-*.tsv* file, see Table 3.

Complete Genomics LFR

Genomic assembly. Genomic assembly was performed as described above for Regular Complete Genomics WGS, with the added assembly step of haplotype generation using well information¹². LFR assemblies include only small variants and their associated haplotypes and well counts.

Assembly metrics are summarized in Table 10. Additional data summary information can be found for each genome in the full package in the *ASM/summary-*.tsv* file, see Table 4.

Ion exome sequencing

Sequencing reads with a mean read length of 190 bp were mapped to human genome version hg19. Mean coverage across AmpliSeq Exome target regions is 256x per run, with raw read accuracy at 99%. Additional summary statistics are in Table 1.

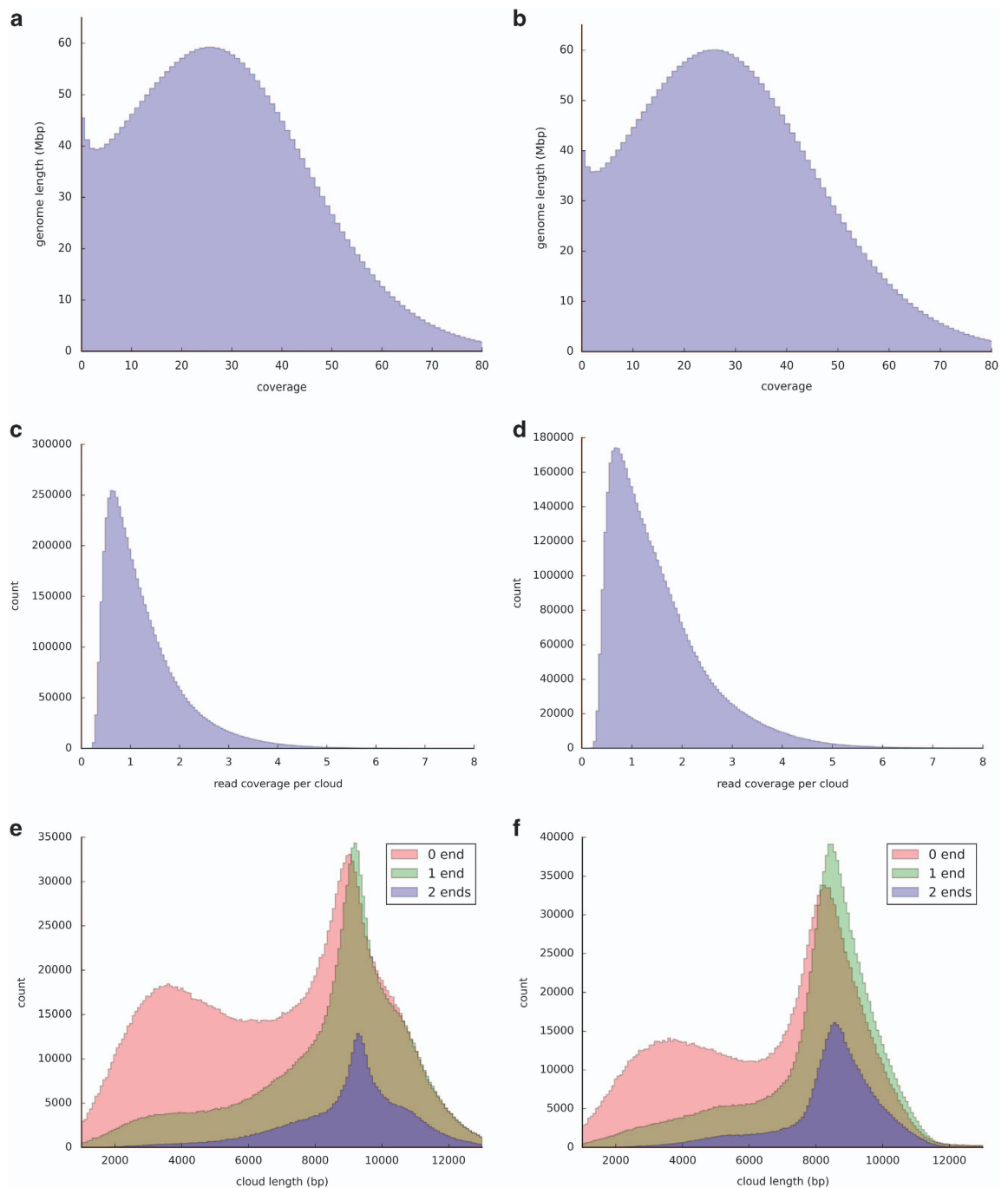


Figure 2. Molecule sequencing characteristics. For HG002 (left—**a,c,e**) and HG003 (right—**b,d,f**), these are distributions of (**a,b**) coverage of the genome by short reads, (**c,d**) read coverage per cloud, and (**e,f**) cloud length for clouds with reads that contain markers at 0, 1, or 2 ends of the cloud.

SOLiD 5500 × I Wildfire WGS

As recommended by the manufacturer, the statistics reported from the instrument from each run of the SOLiD 5500 × I Wildfire WGS were examined to ensure consistency in quality. These statistics included ‘Quality Value’ and ‘Fraction of good+best’, which are a function of the quality of the signal at each ligation. Reads were mapped using Lifescope, which generated the summary statistics given in Table 11.

Bionano Genomics genome maps

De novo assembly of single molecules is accomplished using BioNano Genomics IrysSolve, a proprietary assembler software application, based on an overlap-layout-consensus paradigm^{15–17}. Molecules longer than 150 kb were the input for a pairwise comparison to find all overlaps; then a draft consensus map

	HG002	HG003	HG004	HG005
% Bases in the fastq files which have base quality higher than 30	93.0	93.0	92.9	91.7
Mean target coverage	225	192	213	181
Mean insert size	201	203	200	203
% Duplication	4.9	4.2	4.9	5.7
% Pass Illumina filter reads aligned to the human reference genome	99.7	99.7	99.7	99.8
% Target bases covered by 20 reads or more	98.6	98.4	98.6	98.3
% Target bases covered by 30 reads or more	97.9	97.5	97.9	97.2

Table 7. Illumina paired end WES library statistics. Statistics for each genome were calculated using Picard.

Sample ID	Mean Depth	Median Insert Size	PCR Duplication	Mapping Rate	N50 Linked-Reads per Molecule (LPM)	Molecule Length (mean)	N50 Phase Block	SNPs Phased
HG001	33.9	206	1.5%	96.7%	104	129,913	16,674,432	96.2%
HG002	21.5	205	0.8%	94.5%	71	108,973	12,496,838	98.6%
HG003	25.4	216	1.1%	94.5%	146	165,371	20,501,684	98.8%
HG004	23.8	217	1.2%	93.6%	183	146,362	21,602,191	98.7%

Table 8. 10X Genomics library statistics. Statistics were generated using the GemCode Long Ranger software package.

Subject	NA12878 (HG001)	AJ Son (HG002)	AJ mother (HG004)	AJ father (HG003)	Chinese Son (HG005)	Chinese Son (GM24631)	Chinese mother (GM24695)	Chinese father (GM24694)
Fully called genome fraction	0.975	0.977	0.975	0.976	0.976	0.976	0.974	0.976
Gross mapping yield (Gb)	341	360	376	350	356	353	352	355
Both mates mapped yield (Gb)	313	322	351	313	320	318	314	318
Mate distribution mean	389	395	405	386	395	374	379	387
SNP total count (PASS only)	3449567	3463164	3501307	3436943	3431805	3424251	3445736	3407542
SNP transitions/transversions ratio (PASS only)	2.125	2.127	2.122	2.128	2.126	2.128	2.125	2.127

Table 9. Complete Genomics WGS Summary Metrics. Statistics from Complete Genomics WGS for each genome, including sequencing from candidate NIST RM DNA and from Coriell cells.

Subject	GM12878 (rep1)	GM12878 (rep2)	GM12878 (rep3)	AJ Son (GM24385)	AJ Mother (GM24143)	AJ Father (GM24149)
Fully called genome fraction	0.952	0.968	0.971	0.972	0.968	0.972
Gross mapping yield (Gb)	363	362	368	365	371	363
Both mates mapped yield (Gb)	317	311	316	310	325	312
Mate distribution mean	310	290	290	275	277	282
SNP total count (PASS only)	3308426	3410795	3428944	3443940	3461563	3418049
SNP transitions/transversions ratio (PASS only)	2.134	2.134	2.134	2.135	2.134	2.140
N50	181 kb	257 kb	176 kb	887 kb	1661 kb	734 kb
Median Contig Length	76 kb	84 kb	74 kb	133 kb	149 kb	125 kb
Fraction of Phased SNPs	98.7%	99.4%	99.4%	99.7%	99.6%	99.6%

Table 10. Complete Genomics LFR WGS Summary Metrics. Statistics from Complete Genomics LFR WGS for each genome from Coriell cells, including replicates of GM12878.

Average Per Lane	AJ Son	Chinese Son
AvgNumMapped	69.1	72.7
AvgAlignmentLength	57	55.2
AvgBaseQV	32.9	32.9
AvgNumColorMismatches	4.2	4
AverageCoverage	3.3	2.5
TotalCoverage	78.7	62.2
TotalLanes	24	24

Table 11. SOLiD 5500xl Wildfire WGS statistics. Statistics calculated from Lifescope mapped reads.

	AJ Son (GM24385)	AJ Mother (GM24143)	AJ Father (GM24149)	Chinese Son (GM24631)
Single-molecule maps*				
No. of DNA molecules (k)	1,320	1,107	934	596
Single molecule N50 >150 kb (kb)	264	257	306	314
Average size (kb)	262	257	289	296
Maximum size (kb)	2,930	2,928	2,445	1,858
Total molecule length (Gb)	346	284	270	176
Est. average depth of coverage (X)	112X	92X	87X	57X
Consensus maps (Assembly)				
Total consensus map size (Gb)	2.88	2.86	2.88	2.71
No. consensus maps	1,098	1,092	1,094	3,260
N50 (Mb)	4.59	4.44	4.48	1.09
Longest consensus map size (Mb)	17.7	17.0	25.4	5.9
% genome map aligned to hg19	96%	98%	97%	99%
% Hg19 genome coverage	90%	90%	91%	86%

Table 12. Bionano genome map statistics. Statistics generated from mapping of single molecules as well as the assembly ‘consensus maps’. *Statistics are based on DNA molecules that are greater than 150 kb.

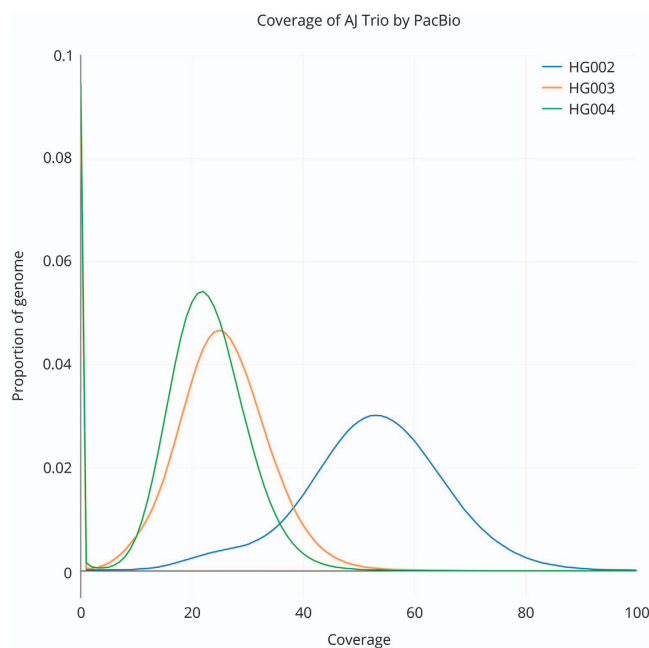


Figure 3. PacBio coverage for AJ Trio. Histogram of coverage from PacBio for the AJ trio generated using bedtools genomecov (raw data available at <https://plot.ly/~justinzook/122/coverage-of-aj-trio-by-pacbio/>).

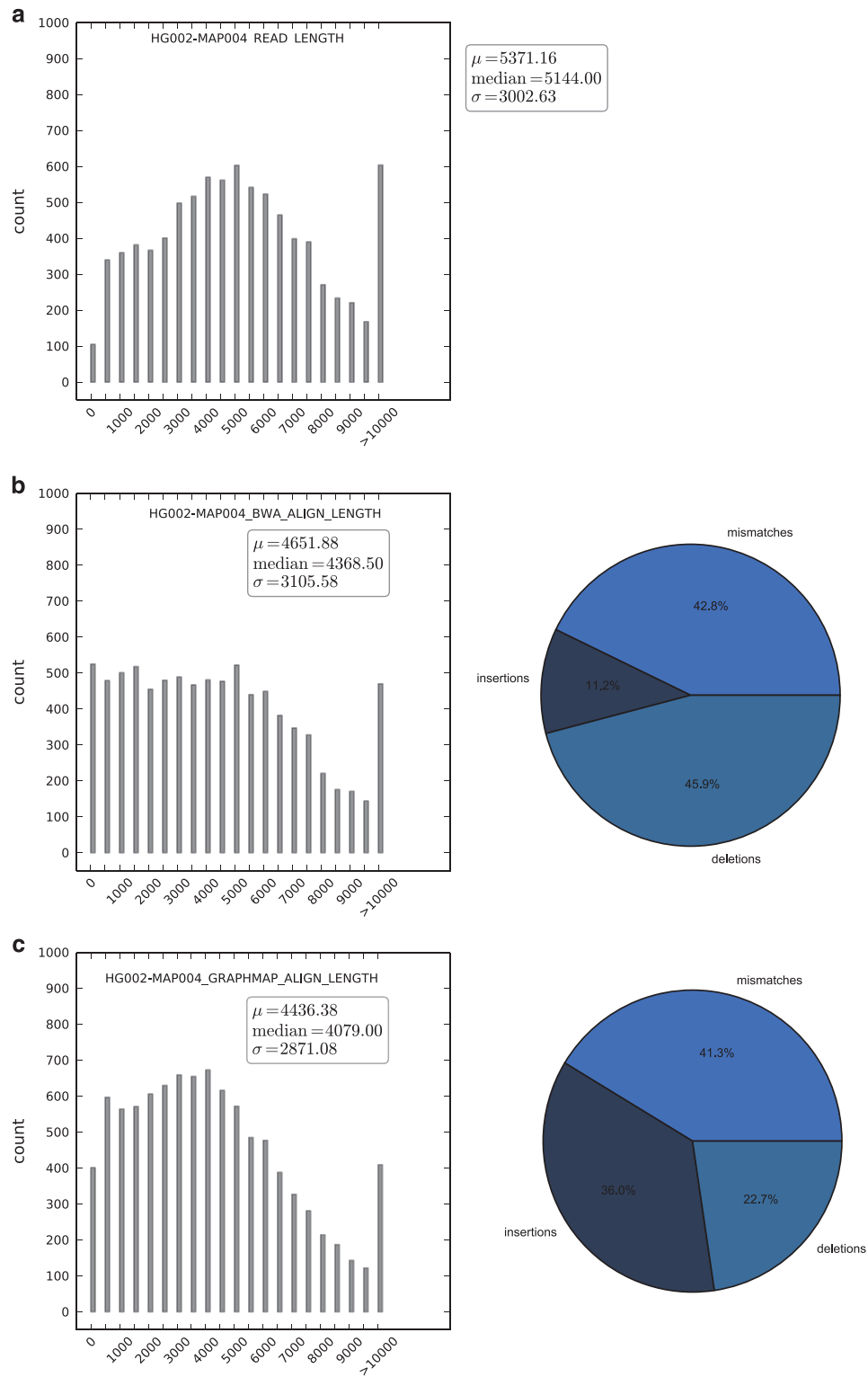


Figure 4. Statistics from first Oxford Nanopore run of the AJ Son (HG002). SQK-MAP-004 sequenced read length distribution (a), then alignment length and alignment error type for BWA (b) and GraphMap (c).

(BioNano Genomics CMAP) was constructed based on these overlaps. The draft BioNano Genomics CMAP was refined by mapping single molecules to it and iteratively recalculating the label positions. Next, the draft BioNano Genomics CMAP (consensus genome maps) were extended by aligning overhanging molecules to the consensus maps and calculating a consensus in the extended regions. Finally, the consensus maps were compared and merged iteratively five times where the patterns matched

	Reads	Bases
	Sequenced	
MAP-006	9677	38651334
MAP-004	9946	46019375
	Mapped with bwa	
MAP-006	9479	35816815
MAP-004	8522	38062962
	Mapped with graphmap	
MAP-006	9625	34474291
MAP-004	9577	40573644

Table 13. Oxford Nanopore read statistics for the AJ Son (HG002). Summary of read count and bases sequenced, as well as mapped with either BWA or GraphMap.

and then the final label position calculation was made. Summary statistics for each sample are presented in Table 12.

PacBio

Assuming a 3.2 Gb human genome, sequencing was conducted to approximately 69X, 32X, and 30X coverage for AJ son (HG002), AJ father (HG003), and AJ mother (HG004) across 292, 139, and 132 SMRT cells, respectively. 27.4, 13.2, and 12.4 M subreads were generated resulting in 220.0, 101.6, and 94.9 Gb of sequence data with sub-read length N50 values of 11,087, 10,728, and 10,629 basepairs. The coverage distribution for each genome is also depicted in Fig. 3.

Oxford Nanopore

The base-calling tool provided by Oxford Nanopore Technologies, Metrichor, is a cloud-based service that provides the user with a choice of 1D or 2D base-calling, the latter being more accurate. Raw current was processed during sequencing of the SQK-MAP-004 library via Metrichor version 2.26 protocol '2D Basecalling' and during sequencing of SQK-MAP-006 library via Metrichor version 2.34.3 protocol '2D Basecalling for SQK-MAP-006.' These protocols segregate output into 'pass' and 'fail' directories, corresponding to the success of alignment of complement and template sections. Basecalls and relevant information are stored in fast5 files in each directory. Poretools (<https://github.com/arq5x/poretools>) was used to extract sequence information from these directories as fasta or fastq files (poretools fasta --type 2D./). The ensuing 2D reads were mapped to the human hg19 reference using BWA (version 0.7.12, parameters: bwa mem -x ont2d) and GraphMap (version 0.22, parameters: default/semiglobal bit-vector alignment mode). CIGAR strings in the generated BAM files were updated to the X/= format using the SamFixCigar module of the jvarkit package (<https://github.com/lindenb/jvarkit/wiki/SamFixCigar>). Error rates were calculated according to the GIGAR strings using a modified version of count-errors.py (<https://github.com/arq5x/nanopore-scripts>), modified code provided as Supplementary Material.

The SQK-MAP-004 and SQK-MAP-006 sequencing runs yielded 46 and 38 Mb of sequence, in 9,946 and 9,677 reads, respectively. Of note is that use of the -type 2D --high-quality option to extract fasta or fastq files from the directories of fast5 files produced by Metrichor may significantly reduce the amount of data available for analysis compared with the data available via only using the --type 2D option. Depending upon application and amount of available data, focusing upon the most high quality subset of 2D reads may be desired, but we have chosen to present all of the 2D reads available from these two sequencing runs here. The SQK-MAP-004 version yielded 2D reads at a median length of 5,144 bp (Fig. 4). GraphMap maps more reads than BWA MEM ont2d, which cover a greater region of the genome (Table 13), and allows a higher error rate (Table 14). The proportion of mismatches to indels was similar for both aligners, however BWA annotates indels preferentially as deletions, and GraphMap favors insertions. The SQK-MAP-006 chemistry yielded reads of median length 3,256 bp (Fig. 5). Again, GraphMap aligns more reads than BWA, though the median alignment length was smaller for GraphMap (Fig. 4) resulting in less total coverage (Table 13). Again, mismatch rates were similar, BWA favors calling deletions, and GraphMap favors calling insertions.

In both cases, GraphMap maps with higher error rate, though, in the case of the the SQK-MAP-006 dataset this does not result in higher coverage. This highlights the importance of the choice of aligner and that overall quality of the reads and nature of analysis may be useful to consider prior to choosing one over the other. While the coverage of these datasets are very low (2X coverage (at least 650 Kb for either dataset, with either mapper) (Table 14).

Usage Notes

The genomes sequenced in this work (Table 2) and their data (Table 1) are all publicly available both as cell lines and as DNA. The pilot genome, NIST RM 8398 (based on Coriell DNA NA12878), is available

	BWA		GRAPHMAP	
	MAP-004	MAP-006	MAP-004	MAP-006
error rate bin (%)	read count by error rate bin			
0	14	374	2	136
10	1975	7314	945	4137
20	4147	1416	2277	1829
30	2160	245	1950	1478
40	226	130	2812	1052
50	0	0	1212	729
60	0	0	344	237
70	0	0	35	27
80	0	0	0	0
90	0	0	0	0
total mapped reads	8522	9479	9577	9625
depth	#bases at depth			
1	37385018	35169774	39664522	33795271
2	475745	446938	612561	427011
3	64330	71522	105609	84999
4	37049	34671	63263	63223
5	23206	21208	40549	30885
6	11790	17722	25127	20379
7	13281	8292	13320	11194
8	10757	5598	14185	10521
9	7136	7936	7809	4821
10	2284	2792	5897	5117
11	2981	3567	4329	4595
12	4739	3038	2766	3925
13	3858	2976	1270	1578
14	2537	3557	1053	1404
15	2117	2850	526	654
16	860	3848	1014	776
17	4148	1496	477	486
18	2002	2258	2510	1235
19	2959	798	945	692
20	3005	507	1242	208
21	1245	894	518	613
22	1273	1260	1720	480
23	530	510	1142	320
24	78	178	1290	728
25	34	137	0	960
26	0	48	0	603
27	0	200	0	388
28	0	149	0	481
29	0	187	0	402
30	0	137	0	57
31	0	110	0	285
32	0	154	0	0
33	0	309	0	0
34	0	204	0	0
35	0	265	0	0
36	0	360	0	0
37	0	300	0	0
38	0	65	0	0

Table 14. Oxford Nanopore mapping and error profile comparison. Error profile and coverage of SQK-MAP-006 and SQK-MAP-004 datasets mapped by either BWA or GraphMap.

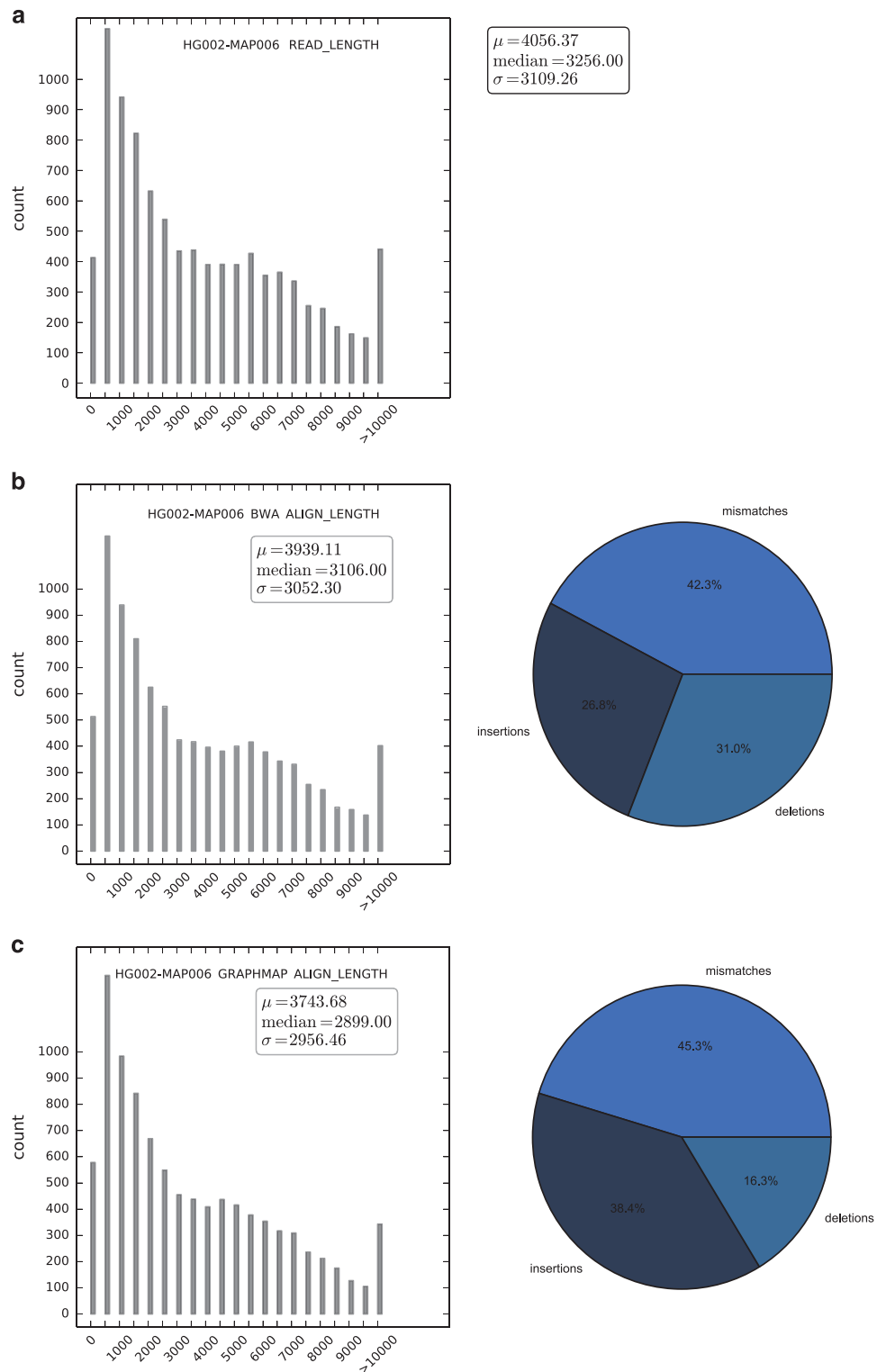


Figure 5. Statistics from second Oxford Nanopore run of the AJ Son (HG002). SK-MAP-006 sequenced read length distribution (a), then alignment length and alignment error type for BWA (b) and GraphMap (c).

both from Coriell as well as from NIST (<http://tinyurl.com/giabpilot>). The NIST RM 8398 was prepared by Coriell from a large growth of cells, and the DNA was extracted and mixed to produce about 8300 10 ug vials of DNA. The remaining genomes are from the Personal Genome Project. These genomes are also available as EBV-immortalized B lymphoblastoid cell lines and as extracted DNA from Coriell, and 4 of them will be available as NIST RMs, planned for release in mid 2016. The AJ Son will be distributed as RM 8391, the AJ Trio will be distributed as RM 8392, and the Chinese Son will be distributed as RM 8393

(note that the Chinese parents are only available from Coriell). Similar to the pilot genome, the other candidate NIST RMs are extracted DNA from a large batch of cells. Except for technologies that optimally start with cells (Complete Genomics LFR, 10X Genomics, and BioNano), all data in this work are collected from the NIST RM DNA. It is possible that small differences may exist between the NIST RM DNA and the DNA from Coriell because they come from different passages of cells and may contain different new mutations.

All data from Genome in a Bottle project are available without embargo, and the primary location for data access is <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp>. To facilitate data analysis in cloud, all the data have been mirrored to the Amazon Web Services Public Datasets repository with 's3://giab' as bucket name. In addition, data that were submitted to SRA can also be accessed through NCBI BioProject (<http://www.ncbi.nlm.nih.gov/bioproject/200694>). The Genome in a Bottle Consortium has formed an Analysis Team to coordinate analyses by groups that are interested in analyzing these data. The primary goal of this group is to establish high-confidence phased variant calls of all sizes for these genomes, so that anyone can benchmark accuracy of their calls for these genomes. The Analysis Team has several sub-groups working on assembly, small variant calling, structural variant calling, and phasing. The intermediate analysis results from these sub-groups are being organized in subdirectory under 'analysis' (<ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/>, <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/analysis/> and <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/analysis/>) with the name describing analyzer's name who performed the analysis, technology for dataset(s) that has been used, type of variant being characterized, analysis tool or algorithm being utilized, and the submission date (MMDDYYYY format) serving as version for better understanding what the datasets were about. The integrated high-confidence calls for the trio samples will be available at <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/>, and the subdirectory with name 'latest' will always contain the latest results published by the Genome in a Bottle Consortium. To make it easier to visualize these data in a web-based genome browser, the GeT-RM browser at NCBI is also hosting some of the vcf and bam files from NA12878 and is currently working on hosting additional data from NA12878 as well as the other GIAB genomes.

In order to improve the accessibility and usability of the Genome in a Bottle project data, a GitHub site (<https://github.com/genome-in-a-bottle>) has also been developed for the Genome in a Bottle project. The idea behind this site is to help the users navigate the main ftp site more easily, and find as much answers as possible by themselves related to the Genome in a Bottle project and its data. Several repositories have been established within the genome-in-a-bottle GitHub site. The repository of 'about_GIAB' describes the objectives and overview of the Genome in a Bottle project, while the repository of 'giab_data_indexes' lists and organizes the index files for the raw sequences and/or alignments by the sequencing platforms for each of the individuals or trio families, therefore, users could use the specific index file to guide their downloading for the desired dataset from the ftp site. The repositories of 'giab_data_analysis' and 'giab_latest_release' provide easy links to the specific ftp locations regarding ongoing analysis results performed by individual analysis groups and the latest high-confidence sets released by the Genome in a Bottle Consortium. The tools and methods that have been used by the analysis groups have been documented in 'giab_tools_methods' repository, while the scientific publications are listed in 'giab_publications' repository. The 'giab_FAQ' repository provides short answers for the frequently asked questions regarding how to effectively access and appropriately utilize the data from the Genome in a Bottle project.

References

1. Rasberry, S. D. & Gills, T. E. The certification, development and use of standard reference materials. *Spectrochim. Acta Part B At. Spectrosc.* **46**, 1577–1582 (1991).
2. Mackey, E. A. *et al.* Certification of NIST Standard Reference Material 1575a Pine Needles and Results of an International Laboratory Comparison. *NIST Special Publication* 260–156 (2004).
3. Lettieri, T. R., Hartman, A. W., Hembree, G. G. & Marx, E. J. Certification of SRM 1960—Nominal 10 micrometer diameter polystyrene spheres (space beads). *Res. Natl. Inst. Stand. Technol.* **96**, 669 (1991).
4. 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).
5. Parikh, H. *et al.* svclassify: a method to establish benchmark structural variant calls. *BMC Genomics* **17**, 64 (2016).
6. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv* <http://arxiv.org/abs/1303.3997> (2013).
7. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
8. Van der Auwera, G. A. *et al.* Current Protocols in Bioinformatics. *Curr. Protoc. Bioinforma* **11**, 11.10.1–11.10.33 (2013).
9. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
10. Drmanac, R. *et al.* Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays. *Science* **327**, 78–81 (2010).
11. Blanco, L. *et al.* Highly efficient DNA synthesis by the phage phi 29 DNA polymerase. Symmetrical mode of DNA replication. *J. Biol. Chem.* **264**, 8935–8940 (1989).
12. Peters, B. A. *et al.* Accurate whole-genome sequencing and haplotyping from 10 to 20 human cells. *Nature* **487**, 190–195 (2012).
13. Carnevali, P. *et al.* Computational Techniques for Human Genome Resequencing Using Mated Gapped Reads. *J. Comp. Bio.* **19**, 279–272 (2012).
14. Faust, G. G. & Hall, I. M. SAMBLASTER: fast duplicate marking and structural variant read extraction. *Bioinformatics* **30**, 2503–2505 (2014).

15. Cao, H. *et al.* Rapid detection of structural variation in a human genome using nanochannel-based genome mapping technology. *Gigascience* **3**, 34 (2014).
16. Valouev, A., Schwartz, D. C., Zhou, S. & Waterman, M. S. An algorithm for assembly of ordered restriction maps from single DNA molecules. *Proc. Natl. Acad. Sci. U. S. A.* **103**, 15770–15775 (2006).
17. Genomic mapping: a statistical and algorithmic analysis of the optical mapping system. *University of Southern California Dissertations and Theses* <http://digitallibrary.usc.edu/cdm/ref/collection/p15799coll127/id/289932> (2010).

Data Citations

1. Zook, J. M. *et al.* NCBI SRA SRX1049768–SRX1049855 (2015).
2. Zook, J. M. *et al.* NCBI SRA SRX847862–SRX848317 (2015).
3. Zook, J. M. *et al.* NCBI SRA SRX1388368–SRX1388459 (2015).
4. Zook, J. M. *et al.* NCBI SRA SRX1388732–SRX1388743 (2015).
5. Sheng, Y. *et al.* NCBI SRA SRP047086 (2015).
6. Schnall-Levin, M. *et al.* NCBI SRA SRX1392293–SRX1392296 (2015).
7. Truty, R. *et al.* NCBI SRA SRX840234 (2014).
8. Truty, R. *et al.* NCBI SRA SRX852932–SRX852936 (2014).
9. Hyland, F. *et al.* NCBI SRA SRX847094 (2014).
10. Hyland, F. *et al.* NCBI SRA SRX848742–SRX848744 (2014).
11. Hyland, F. *et al.* NCBI SRA SRX326642 (2013).
12. Zook, J. M. *et al.* NCBI SRA SRX1497273 (2015).
13. Zook, J. M. *et al.* NCBI SRA SRX1497276 (2015).
14. Sebra, R. *et al.* NCBI SRA SRX1033793–SRX1033798 (2015).

Acknowledgements

We thank members of the Genome in a Bottle Consortium for their feedback before and during this study. C.E.M. acknowledges funding from the Starr Cancer Consortium grant I4-A442, I7-A765, I9-A9-071), the Irma T. Hirsch and Monique Weill-Caulier Charitable Trusts, Bert L and N Kuggie Vallee Foundation and the WorldQuant Foundation (CEM), The Pershing Square Sohn Cancer Research Alliance, NASA (NNX14AH50G), and the National Institutes of Health (R25EB020393, R01NS076465). Certain commercial equipment, instruments, or materials are identified in this paper only to specify the experimental procedure adequately. Such identification is not intended to imply recommendation or endorsement by the NIST, nor is it intended to imply that the materials or equipment identified are necessarily the best available for the purpose.

Author Contributions

J.M.Z., M.L.S., and G.I.A.B. designed the overall study. J.M.Z., D.C., J.M., L.V., N.S., Z.W., Y.L., N.A., E.H., E.J., R.S., R.M.T., K.Z., Y.F., M.C., C.X., and M.L.S. wrote the manuscript. J.M.Z., D.C., J.M., L.V., and M.L.S. designed, sequenced, and analyzed the Illumina paired end W.G.S. J.M.Z., D.C., J.M., N.S., A.S., Z.W., and M.L.S. designed, sequenced, and analyzed the Illumina mate pair W.G.S. J.M.Z., D.C., J.M., N.S., A.S., Y.L., F.C., E.J., A.M., and M.L.S. designed, sequenced, and analyzed the Illumina synthetic long read W.G.S. C.E.M., N.A., and E.H. designed, sequenced, and analyzed the Oxford Nanopore W.G.S. K.P., W.S., T.L., M.S., Z.D., A.H., and H.C. designed, sequenced, and analyzed the BioNano mapping. R.M.T., C.C.C., and N.G. designed, sequenced, and analyzed the Complete Genomics W.G.S. K.Z., S.G., F.H., and Y.F. designed, sequenced, and analyzed the Ion Torrent exome sequencing. J.M.Z., J.M., G.D., E.S., R.S., A.B., M.C., and M.L.S. designed, sequenced, and analyzed the PacBio W.G.S. J.M.Z., A.W.Z., M.B., J.B., P.E., G.M.C., M.L.S., and G.I.A.B. designed the process for selecting the genomes from the P.G.P. P.M., S.K.-P., G.S.Y.Z., M.S.-L., H.S.O., and P.A.M. designed, sequenced, and analyzed the 10X Genomics data. Y.S. and K.B.R. designed, sequenced, and analyzed the Illumina W.E.S. data. J.M.Z., S.S., J.T., and C.X. designed and manage the GIAB FTP site and SRA submissions.

Additional Information

Supplementary information accompanies this paper at <http://www.nature.com/sdata>

Competing financial interests: F.C., E.J., A.M. are employees of Illumina. K.P., W.S., T.L., M.S., Z.D., A.H., and H.C. are employees of BioNano Genomics. P.M., S.K.-P., G.S.Y.Z., M.S.-L., H.S.O., and P.A.M. are employees of 10X Genomics. R.M.T., C.C.C., and N.G. are employees of BGI-Complete Genomics. K.Z., S.G., F.H., and Y.F. are employees of Thermo Fisher Scientific.

How to cite this article: Zook, J. M. *et al.* Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**:160025 doi: 10.1038/sdata.2016.25 (2016).



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

Metadata associated with this Data Descriptor is available at <http://www.nature.com/sdata/> and is released under the CC0 waiver to maximize reuse.