

An open resource for accurately benchmarking small variant and reference calls

Justin M. Zook^{1*}, Jennifer McDaniel¹, Nathan D. Olson¹, Justin Wagner¹, Hemang Parikh¹, Haynes Heaton^{2,3}, Sean A. Irvine⁴, Len Trigg⁴, Rebecca Truty⁵, Cory Y. McLean^{6,7}, Francisco M. De La Vega⁸, Chunlin Xiao⁹, Stephen Sherry⁹ and Marc Salit^{1,10,11}

Benchmark small variant calls are required for developing, optimizing and assessing the performance of sequencing and bioinformatics methods. Here, as part of the Genome in a Bottle (GIAB) Consortium, we apply a reproducible, cloud-based pipeline to integrate multiple short- and linked-read sequencing datasets and provide benchmark calls for human genomes. We generate benchmark calls for one previously analyzed GIAB sample, as well as six genomes from the Personal Genome Project. These new genomes have broad, open consent, making this a ‘first of its kind’ resource that is available to the community for multiple downstream applications. We produce 17% more benchmark single nucleotide variations, 176% more indels and 12% larger benchmark regions than previously published GIAB benchmarks. We demonstrate that this benchmark reliably identifies errors in existing callsets and highlight challenges in interpreting performance metrics when using benchmarks that are not perfect or comprehensive. Finally, we identify strengths and weaknesses of callsets by stratifying performance according to variant type and genome context.

Genome sequencing is increasingly used in clinical applications, making variant calling accuracy of paramount importance. With this in mind, the GIAB developed benchmark small variants for the pilot GIAB genome¹. These benchmark calls have been used for optimization and analytical validation of clinical sequencing^{2–4}, comparisons of bioinformatics tools⁵ and the optimization, development and demonstration of new technologies⁶.

Here, we build on previous GIAB integration methods to enable the development of highly accurate and reproducible benchmark genotype calls from any genome, using multiple datasets from different sequencing methods (Fig. 1). We first generate more comprehensive and accurate integrated single nucleotide variations (SNVs), small indel and homozygous reference calls for HG001—the same sample used in the original GIAB analysis. These calls were made more accurate and comprehensive by the use of new and higher-coverage short- and linked-read datasets⁷, technology-optimized variant calling methods and more robust methods to determine when to trust a variant call or homozygous reference region from each technology. Calls supported by two technologies were used to train a model that identified calls from each dataset that were outliers and therefore less trustworthy. Variant calls and regions were included in the benchmark set if at least one input callset was trustworthy and all trustworthy callsets agreed. We compare our callsets to the phased pedigree-based callsets from the Illumina Platinum Genomes Project⁸ (PGP) and ref. ⁹, following best practices established by the Global Alliance for Genomics and Health Benchmarking Team¹⁰, and manually examine the differences to evaluate the accuracy of our callsets.

We apply this approach to six broadly consented GIAB genomes from the PGP¹¹, an Ashkenazi Jewish mother–father–son trio and a Han Chinese mother–father–son trio. We also use our methods to form similar benchmark sets with respect to the human GRCh38 reference.

We demonstrate that our benchmark sets reliably identify false-positive and -negative variant calls in a variety of high-quality variant callsets, including alternative benchmark sets such as Platinum Genomes (PG). The benchmark sets can be used to identify strengths and weaknesses of different methods. Finally, we discuss how strengths and weaknesses of the benchmark sets themselves can bias the resulting performance metrics, particularly for difficult variant types and genome contexts.

In contrast to existing benchmarking efforts, the extensively characterized PGP genomes analyzed here have an open, broad consent and broadly available data and cell lines, as well as consent to re-contact for additional types of samples¹¹. They are an enduring resource uniquely suitable for diverse research and commercial applications. In addition to homogeneous DNA from National Institute of Science and Technology Reference Materials (NIST RMs) and DNA and cell lines from Coriell, a variety of products using these genomes is already available, including induced pluripotent stem cells, cell line mixtures, cell line DNA with synthetic DNA spike-ins with mutations of clinical interest, formalin-fixed paraffin-embedded cells and circulating tumor DNA mimics^{12,13}.

Previous studies have used more restricted samples to characterize variants and regions complementary to our benchmark sets: for example, the PG pedigree analysis⁸, HuRef analysis using Sanger sequencing^{14,15}, integration of multiple structural variant (SV) calling

¹Material Measurement Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, USA. ²10x Genomics, Pleasanton, CA, USA.

³Wellcome Trust Sanger Institute, Hinxton, Cambridge, UK. ⁴Real Time Genomics, Hamilton, New Zealand. ⁵Invitae Corporation, San Francisco, CA, USA.

⁶Verily Life Sciences, South San Francisco, CA, USA. ⁷Google Inc., Mountain View, CA, USA. ⁸Department of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA. ⁹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA. ¹⁰Joint Initiative for Metrology in Biology, Stanford, CA, USA. ¹¹Department of Bioengineering, Stanford University, Stanford, CA, USA.

*e-mail: jzook@nist.gov

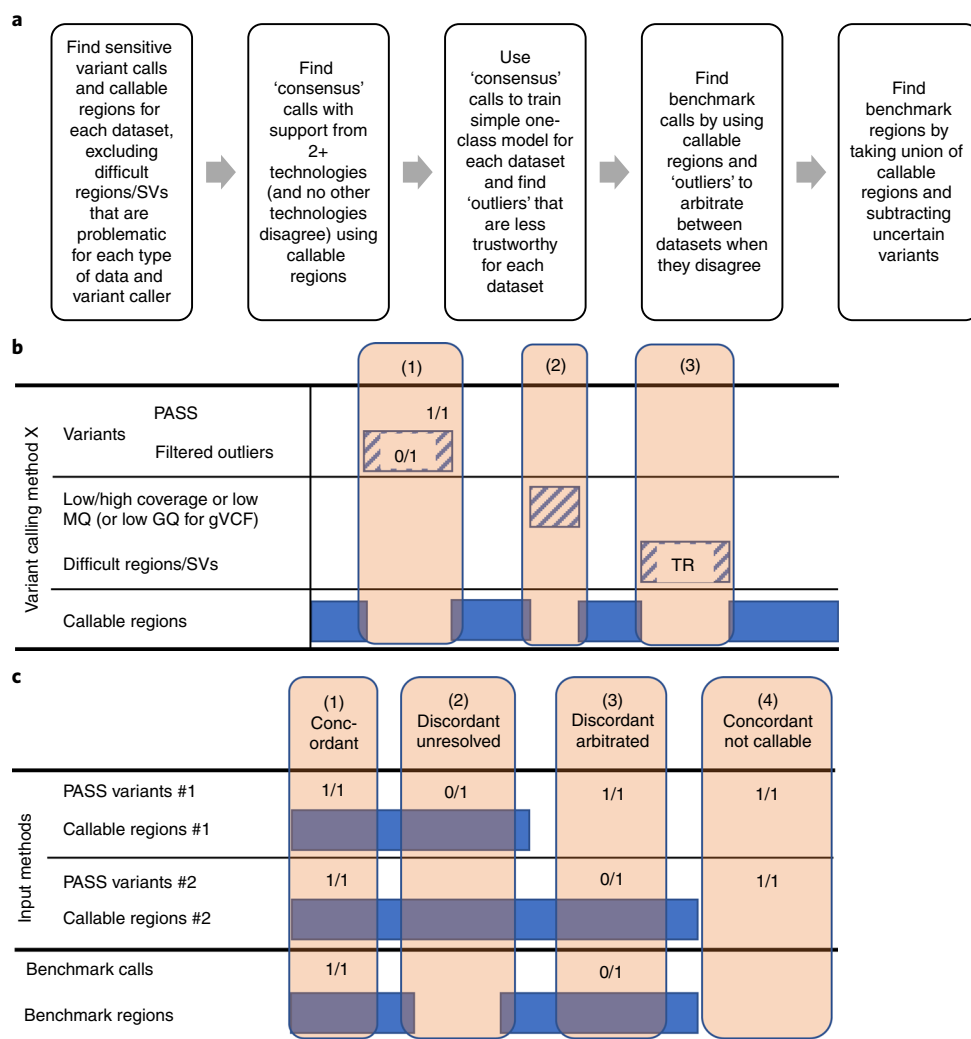


Fig. 1 | Arbitration process used to form our benchmark set from multiple technologies and callsets. **a**, The arbitration process has two cycles. The first cycle ignores ‘filtered outliers’. Calls that are supported by at least two technologies in the first cycle are used to train a model that identifies variants from each callset with any annotation value that is an ‘outlier’ compared to these two-technology calls. In the second cycle, the outlier variants and surrounding 50 base pairs (bp) are excluded from the callable regions for that callset. **b**, For each variant calling method, we delineate callable regions by subtracting regions around filtered/outlier variants as at locus (1), regions with low coverage or mapping quality (MQ) as at locus (2) and ‘difficult regions’ prone to systematic miscalling or missing variants for the particular method as at locus (3). For callsets in genetic variant call format (gVCF), we exclude homozygous reference regions and variants with genotype quality (GQ) < 60. Difficult regions include different categories of tandem repeats (TR) and segmental duplications. **c**, Four arbitration examples with two arbitrary input methods. (1) Both methods have the same genotype and variant and this is in their callable regions, so the variant and region are included in the benchmark set. (2) Method 1 calls a heterozygous variant and Method 2 implies homozygous reference, and it is in both methods’ callable regions, so the discordant variant is not included in the benchmark calls and 50 bp on each side are excluded from the benchmark regions. (3) The methods have discordant genotypes, but the site is only inside Method 2’s callable regions, so the heterozygous genotype from Method 2 is trusted and is included in the benchmark regions. (4) The two methods’ calls are identical, but because they are outside both methods’ callable regions the site is excluded from benchmark variants and regions.

methods on HS1011 (ref. ¹⁶) and synthetic diploid using long-read sequencing of hydatiform moles¹⁷. The recent synthetic diploid method¹⁷ has the advantage that long-read assemblies of each haploid cell line have orthogonal error profiles to the diploid short- and linked-read-based methods used in this manuscript, though short reads were used to correct indel errors in the long-read assemblies. However, the synthetic diploid hydatiform mole cell lines are not available in a public repository¹⁰. The broadly consented and available GIAB samples from PGP have at least three strengths relative to the other samples: (1) we can continue to sequence these renewable cell lines with new technologies and improve our benchmark set over time; (2) any clinical or research laboratory can sequence these

samples with their own method and compare their results to our benchmark set; and (3) a wide array of secondary reference samples can be generated from the same cell lines to meet particular community needs.

Results

Design of benchmark calls and regions. Our goal is to design a reproducible, robust and flexible method to produce a benchmark variant and genotype callset (including homozygous reference regions). When any sequencing-based callset is compared to our benchmark calls requiring stringent matching of alleles and genotypes, the majority of discordant calls in the benchmark regions

Table 1 | Summary of statistics of GIAB benchmark calls and regions from HG001 from v.2.18 to v.3.3.2 and their comparison to Illumina PG 2016-v.1.0 calls

Integration version	v.2.18	v.2.19	v.3.2	v.3.2.2	v.3.3	v.3.3.1	v.3.3.2	v.3.3.1	v.3.3.2
Reference	GRCh37	GRCh37	GRCh37	GRCh37	GRCh37	GRCh37	GRCh37	GRCh38	GRCh38
Integration date	Sep 2014	Apr 2015	May 2016	June 2016	Aug 2016	Oct 2016	Nov 2016	Oct 2016	Nov 2016
Number of bases in benchmark regions (chromosomes 1–22 + X; Gb)	2.20	2.22	2.54	2.53	2.57	2.58	2.58	2.45	2.44
Fraction of non-N bases covered in chromosomes 1–22 + X (%)	77.4	78.1	89.6	89.2	90.5	91.1	90.8	84.2%	83.8
Fraction of RefSeq coding sequence covered (%)	73.9	74.0	87.8	87.9	89.9	90.0	89.9	83.4	83.3
Total number of calls in benchmark regions	2,915,731	3,153,247	3,433,656	3,512,990	3,566,076	3,746,191	3,691,156	3,617,168	3,542,487
Single-nucleotide polymorphisms	2,741,359	2,787,291	3,084,406	3,154,259	3,191,811	3,221,456	3,209,315	3,058,368	3,042,789
Insertions	86,204	172,671	171,866	176,511	171,715	243,856	225,097	269,331	241,176
Deletions	87,161	189,932	169,389	173,976	189,807	266,386	245,552	275,041	247,178
Block substitutions	1,005	2,532	7,476	7,716	10,364	13,332	11,192	13,976	11,344
Transition/transversion ratio for SNVs	2.12	2.12	2.14	2.14	2.11	2.10	2.10	2.10	2.11
Fraction phased globally or locally (%)	0.0	0.3	3.9	3.9	8.8	99.0	99.6	98.5	99.5
Comparisons of GIAB to PG calls									
Number of GIAB calls concordant with PG in both PG and GIAB beds	2,825,803	3,030,703	3,312,580	3,391,783	3,441,361	3,550,914	3,529,641	3,459,674	3,431,752
Number of PG-only calls in both beds	194	404	81	52	60	67	61	202	180
Number of GIAB-only calls in both beds	49	87	56	57	40	50	47	105	94
Number of PG-only calls	1,223,697	1,018,795	274,671	138,894	550,982	445,563	469,202	659,870	690,887
Total after excluding PG bed	605,142	493,694	118,810	22,224	172,375	140,857	142,682	386,739	391,523
Number of GIAB-only calls in GIAB benchmark bed	90,722	122,544	121,076	121,207	124,715	195,277	163,467	157,494	111,787
Number of concordant calls filtered by the GIAB benchmark bed	12	0	736,918	657,715	608,137	53,460	51,255	48,696	45,779

Note that PG bed files were contracted by 50 bp to minimize partial complex variant calls in the PG calls.

(that is, false-positives and -negatives) should be attributable to errors in the sequencing-based callset. We develop a modular, cloud-based data integration pipeline, enabling diverse data types to be integrated for each genome. We produce benchmark variant calls and regions by integrating methods and technologies that have different strengths and limitations, using evidence of potential bias to arbitrate when methods have differing results (Fig. 1). Finally, we evaluate the utility of the benchmark variants and regions by comparing high-quality callsets to the benchmark calls and manually curating discordant calls to ensure that most are errors in the other high-quality callsets.

New benchmark sets are more comprehensive and accurate. The new integration method is designed to give more accurate and comprehensive benchmark sets in several ways, which are detailed at the end of Methods: (1) We use newer input data, including 10x Genomics linked reads, and most data were measured from the homogeneous NIST RM batches of DNA. (2) Some long homopolymers and difficult-to-map regions are now in the benchmark sets by using the PCR-free Illumina Genome Analysis Toolkit (GATK) gVCF outputs and 10x Genomics LongRanger, respectively. (3)

We use a more accurate set of potential structural variants to exclude regions with problematic small variant calls. (4) We add phasing information. Table 1 shows the evolution of the GIAB/NIST benchmark sets since our previous publication¹. The fraction of non-N bases in the GRCh37 reference covered has increased from 77.4 to 90.8%, and the numbers of benchmark SNVs and indels have increased by 17 and 176%, respectively. The fraction of GRCh37 Reference Sequence (RefSeq) coding regions covered has increased from 73.9 to 89.9%. The gains in benchmark regions and calls are a result of using new datasets and a more accurate and less conservative SV bed file (Methods), as well as a new approach where excluded difficult regions are callset-specific depending on read length, error profiles and analysis methods. GRCh38 is currently characterized less comprehensively than GRCh37 (Supplementary Table 1, Supplementary Fig. 1 and Supplementary Note 1). To determine whether fewer input data could be used, we also integrated only two datasets (Illumina and 10x), resulting in a larger number of benchmark indel calls but a fivefold higher error rate compared to v.3.3.2 (Supplementary Note 2). Our callsets are now mostly phased using pedigree or trio analysis, but we did not evaluate the accuracy of this phasing (Supplementary Note 3).

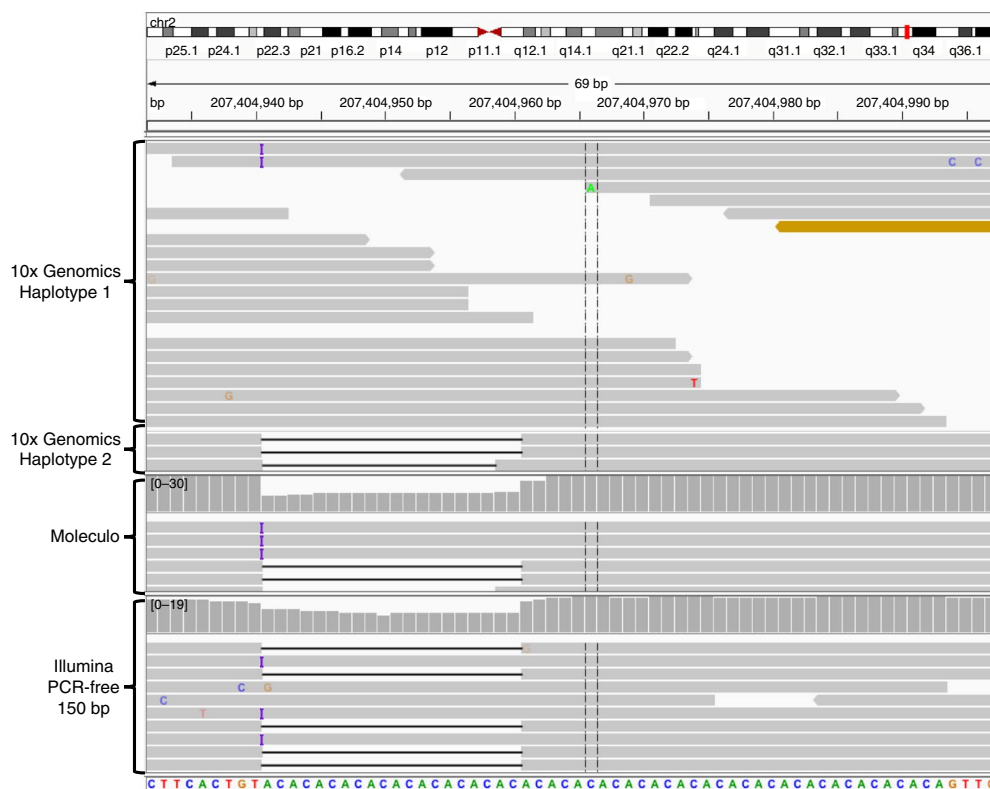


Fig. 2 | Complex variant discordant between GIAB and Illumina PG. Compound heterozygous insertion and deletion in HG001 in a tandem repeat at 2:207404940 (GRCh37), for which Illumina PG calls only a heterozygous deletion. When a callset with the true compound heterozygous variant is compared to PG, it is counted as both a false-positive and a false-negative. Both the insertion and deletion are supported by PCR-free Illumina (bottom) and Molecu-assembled long reads (middle), while reads assigned haplotype1 in 10x support the insertion and reads assigned haplotype2 in 10x support the deletion (top).

High concordance with Illumina PG. The lower section of Table 1 shows increasing concordance over multiple versions of our HG001 callset with the Illumina PG 2016-v.1.0 callset⁸. PG provides orthogonal confirmation of our calls, since the PG phased pedigree analysis may identify different biases from our method. PG contains a larger number of benchmark variant calls, even relative to our newest v.3.3.2 calls for HG001. However, when comparing v.3.3.2 to PG within both benchmark regions, we found that the majority of differences were places where v.3.3.2 was correct and PG had a partial complex or compound heterozygous call (for example, Fig. 2), or where the PG benchmark region partially overlapped with a true deletion that v.3.3.2 called correctly. To reduce the number of these inaccurate calls in PG, we contracted each of the PG benchmark regions by 50 bp on each side, which is similar to how we remove 50 bp on each side of uncertain variants (for example, locus (1) in Fig. 1b). Because PG has many more uncertain variants, contracting the PG benchmark regions by 50 bp reduces the number of variants in the PG benchmark region by 32%, but it also eliminates 93% of the differences between v.3.3.2 and PG in GRCh37, so that 61 PG-only and 47 v.3.3.2-only calls remain in both benchmark regions.

After manually curating the remaining 108 differences between v.3.3.2 and PG in GRCh37, v.3.3.2 has five clear false-negatives, one clear false-positive, two unclear potential false-negatives and six unclear potential false-positives in the NIST benchmark regions (Supplementary Data 1 and Supplementary Note 4). Based on 3,529,641 variants concordant between NIST and PG, there would be about two false-positives and two false-negatives per million true variants. Since these are only in the PG benchmark regions, these error rates are probably lower than the overall error rates, which are difficult to estimate.

Trio Mendelian inheritance analysis. Benchmark callsets were generated for two PGP trios, an Ashkenazi Jewish and a Han Chinese (Supplementary Table 1 and Supplementary Fig. 1). As orthogonal benchmark callsets are not available for these genomes, a trio Mendelian inheritance analysis was used to evaluate the callsets. The trio analysis allows us to phase some variants in the sons (see below) and determine variants that have a genotype pattern inconsistent with Mendelian inheritance.

We separately analyzed Mendelian inconsistent variants that were potential cell line or germline de novo mutations (that is, the son was heterozygous and both parents were homozygous reference), and those that had any other Mendelian inconsistent pattern (which are unlikely to have a biological origin). Out of 2,038 Mendelian violations in the Ashkenazi son, 1,110 SNVs and 213 indels were potential de novo mutations. Out of 425 Mendelian violations in the Chinese son, 103 SNVs and 43 indels were potential de novo mutations. The large difference between the trios may have arisen from differences in the datasets used to generate the benchmark callsets or mutations arising in the Epstein–Barr virus immortalization process. The number of de novo SNVs is about 10% higher than the 1,001 de novo SNVs found by 1000 Genomes for NA12878/HG001, and higher than the 669 de novo SNVs they found for NA19240 (ref. ¹⁸). Following manual inspection of ten random de novo SNVs from each trio, 17/20 appeared to be true de novo. After manual inspection of ten random de novo indels from each trio, 10/20 appeared to be true de novo indels (most in homopolymers or tandem repeats) and the remainder were correctly called in the sons but missed in one of the parents. Supplementary Data 1 and Supplementary Note 5 contain additional details on the manual curation of each site and Mendelian concordance.

Based on the Mendelian analyses, there may be as many as 16SNV and 150indel errors per million variants in one or more of the three individuals before excluding these sites. Note that we exclude the errors found by the Mendelian analysis from our final benchmark bed files, but it is likely that the error modality that caused these Mendelian errors also led to other errors that were not found to be Mendelian inconsistent. Many other error modalities can be Mendelian consistent (for example, systematic errors that cause all genomes in the pedigree to be heterozygous, and systematic errors that cause the wrong allele to be called), so these error estimates should be interpreted with caution.

Performance metrics vary in comparison to different benchmark callsets. We next explore how performance metrics can differ when using different benchmark sets. When comparing a method's callset against these imperfect and incomplete benchmark callsets, estimates of a method's precision and recall/sensitivity computed from imperfect and incomplete benchmark callsets may differ from the true precision and recall of the method for all regions of the genome and all types of variant. Insufficient examples of variants of particular types can result in estimates of precision and recall with high uncertainties. In addition, estimates of precision and recall can be biased by the following: (1) errors in the benchmark callset, (2) use of methods that do not compare differing representations of complex variants and (3) biases in the benchmark towards easier variants and genome contexts. We demonstrate that these biases are particularly apparent for certain challenging variant types and genome contexts, by benchmarking a standard Burrows-Wheeler Aligner-Genome Analysis Tool Kit (bwa-GATK) callset against three benchmark sets: GIABv.2.18 from our previous publication¹, GIABv.3.3.2 from this publication and PG 2016–1.0⁸.

We examine three striking examples in coding regions, complex variants and decoy-associated regions: (1) The bwa-GATK SNV false-negative rate is 62-fold higher when benchmarking PG against GIAB, and 2/10 manually curated false-negatives versus PG appeared to be errors in PG. (2) The bwa-GATK putative false-positive rate for compound heterozygous indels was only slightly lower for PG versus GIAB, but following manual curation 6/10 putative false-positives were actually errors in PG whereas only 1/10 putative false-positives was an error in GIAB. (3) The 1000Genomes developed the hs37d5 decoy sequence to remove false-positives in some regions with segmental duplications not in GRCh37¹⁹. The precision for bwa-GATK-nodecoy is 67% versus GIABv.3.3.2, whereas it is 91% versus GIABv.2.18 and 93% versus PG. The much higher precision versus v.2.18 or PG is a result of many of the decoy-associated variants being excluded from the benchmark regions of GIABv.2.18 and PG. Assessing these variants is important because they comprise 7,123 of the 16,452 false-positives (43%) versus GIABv.3.3.2 for this callset. We also found notable differences between benchmarks for large indels (Supplementary Data 1 and Supplementary Note 6), and subtle differences between performance metrics for different GIAB genomes (Supplementary Note 7 and Supplementary Table 2). Since characteristics of the benchmark can affect performance metrics, we have outlined known limitations of our benchmark set (Supplementary Note 8).

Orthogonal 'validation' of benchmark set. We have chosen to present detailed results of expert manual curation of multiple short-, linked- and long-read sequencing alignments in place of the more traditional targeted Sanger sequencing method for validation. There are several reasons why confirming some discrepancies by Sanger sequencing is unlikely to add more confidence: (1) We already have multiple technologies supporting many of the calls in our truth set. (2) In our manual curation, we visualized PacBio and Moleculo sequencing data, which have longer reads (or pseudo-reads) than Sanger sequencing and were not used to form our truth sets. (3) All 23 of the calls we classified as 'Incorrect or partial call by PG at or near a true variant' are in

homopolymers or tandem repeats, and 20 of these are heterozygous or compound heterozygous indels that Sanger could not help resolve. (4) We classified 18 SNVs as 'Calls that are in PG and appear likely to be FPs despite inheriting properly, either due to a local duplication or a systematic sequencing error that occurs on only one haplotype' (Supplementary Data 1). These regions are unlikely to be easily characterized by Sanger sequencing, and longer reads from PacBio and Moleculo sequencing are more useful orthogonal validation. (5) Even when using a single, short-read sequencing technology, previous work found that discrepancies between Sanger sequencing and their variant calls were mostly errors in Sanger sequencing²⁰. (6) For 32 of the 47 variants found only in v.3.3.2, the Real Time Genomics (RTG) phased pedigree calls had an identical call to v.3.3.2 even though RTG's Illumina-only pedigree methods are much more similar to PG's pedigree methods. For these reasons, we have focused our efforts on expert manual curation of multiple whole-genome sequencing technologies rather than performing targeted Sanger sequencing.

Discussion

Trustworthy benchmarks reliably identify errors in the results of any method being compared to the benchmark. To demonstrate that we meet this goal, we show that when using our benchmark calls with best practices for benchmarking, established by the Global Alliance for Genomics and Health (Supplementary Note 9)¹⁰, the majority of putative false-positives and -negatives are errors in the callset being benchmarked. This is a moving target, since new methods are continually being developed that may perform better, particularly in challenging regions of the genome.

We form our benchmark callsets to accurately represent genomic variation in regions where we are confident in our technical and analytical performance. One notable limitation in such callsets is that they will tend to exclude more difficult types of variation and regions of the genome. When new methods with different characteristics than those used to form our callsets are applied, they may be more accurate in those regions difficult to access with our methods, and less accurate in regions where we are quite confident. For example, as graph-based variant calling methods are developed, these may be able to make much better calls in regions with many alternative locus (ALT) haplotypes such as the major histocompatibility complex, even if they have lower accuracy within our benchmark regions. Therefore, it can be important to consider variant calls outside our benchmark regions when evaluating the accuracy of different variant calling methods.

To assure that our benchmark calls continue to be useful as sequencing and analysis methods improve, we are exploring several approaches for characterizing more difficult regions of the genome. Long reads and linked reads show promise in enabling benchmark calls in difficult-to-map regions of the genome^{21,22}. The methods presented here are a framework to integrate these and other new data types, providing that variant calls with high sensitivity in a specified set of genomic regions can be generated, and that filters could be applied to give high specificity.

Finally, to ensure that these data are an enduring resource that we can improve over time, we have made available an online form where we and others can enter potentially questionable sites in our benchmark regions (<https://goo.gl/forms/zvxjRsYTdrkhqdzM2>). The results of this form are public and updated in real time, so that anyone can see where others have manually reviewed or interrogated the evidence at any site (<https://docs.google.com/spreadsheets/d/1kHgRLinYcnxX3-ulvijf2HrIdrQWz5R5PtXZS-s6ZM/edit?usp=sharing>).

Online content

Any methods, additional references, Nature Research reporting summaries, source data, statements of data availability and associated accession codes are available at <https://doi.org/10.1038/s41587-019-0074-6>.

Received: 25 May 2018; Accepted: 19 February 2019;
Published online: 1 April 2019

References

1. 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).
2. Patwardhan, A. et al. Achieving high-sensitivity for clinical applications using augmented exome sequencing. *Genome Med.* **7**, 71 (2015).
3. Lincoln, S. E. et al. A systematic comparison of traditional and multigene panel testing for hereditary breast and ovarian cancer genes in more than 1000 patients. *J. Mol. Diagnostics* **17**, 533–544 (2015).
4. Telenti, A. et al. Deep sequencing of 10,000 human genomes. *Proc. Natl Acad. Sci. USA* **113**, 11901–11906 (2016).
5. Cornish, A. & Guda, C. A comparison of variant calling pipelines using Genome in a Bottle as reference. *Biomed. Res. Int.* **2015**, 1–11 (2015).
6. Poplin, R. et al. A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983 (2018).
7. Zook, J. M. et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
8. Eberle, M. A. et al. A reference data set of 5.4 million phased human variants validated by genetic inheritance from sequencing a three-generation 17-member pedigree. *Genome Res.* **27**, 157–164 (2017).
9. Cleary, J. G. et al. Joint variant and *de novo* mutation identification on pedigrees from high-throughput sequencing data. *J. Comput. Biol.* **21**, 405–419 (2014).
10. Krusche, P. et al. Best practices for benchmarking germline small variant calls in human genomes. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-019-0054-x> (2019).
11. Ball, M. P. et al. A public resource facilitating clinical use of genomes. *Proc. Natl Acad. Sci. USA* **109**, 11920–11927 (2012).
12. Kudalkar, E. M. et al. Multiplexed reference materials as controls for diagnostic next-generation sequencing: a pilot investigating applications for hypertrophic cardiomyopathy. *J. Mol. Diagn.* **18**, 882–889 (2016).
13. Lincoln, S. E. et al. An interlaboratory study of complex variant detection. Preprint at *bioRxiv* <https://doi.org/10.1101/218529> (2017).
14. Zhou, B. et al. Extensive and deep sequencing of the Venter/HuRef genome for developing and benchmarking genome analysis tools. *Sci. Data* **5**, 180261 (2018).
15. Mu, J. C. et al. Leveraging long read sequencing from a single individual to provide a comprehensive resource for benchmarking variant calling methods. *Sci. Rep.* **5**, 14493 (2015).
16. English, A. C. et al. Assessing structural variation in a personal genome—towards a human reference diploid genome. *BMC Genomics* **16**, 286 (2015).
17. Li, H. et al. A synthetic-diploid benchmark for accurate variant-calling evaluation. *Nat. Methods* **15**, 595–597 (2018).
18. Conrad, D. F. et al. Variation in genome-wide mutation rates within and between human families. *Nat. Genet.* **43**, 712–714 (2011).
19. Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
20. Beck, T. F. et al. Systematic evaluation of Sanger validation of next-generation sequencing variants. *Clin. Chem.* **62**, 647–654 (2016).
21. Marks, P. et al. Resolving the full spectrum of human genome variation using linked-reads. Preprint at *bioRxiv* <https://doi.org/10.1101/230946> (2018).
22. Wenger, A. M. et al. Highly-accurate long-read sequencing improves variant detection and assembly of a human genome. Preprint at *bioRxiv* <https://doi.org/10.1101/519025> (2019).

Acknowledgements

We thank the many contributors to GIAM Consortium discussions. We especially thank R. Saldana and the Sentieon team for advice on running the Sentieon pipeline; A. Carroll and the DNAnexus team for advice on implementing the pipeline in DNAnexus; F. Hyland, S. Ghosh, K. Zhao and J. Bodeau at ThermoFisher for advice on integrating Ion exome and SOLiD genome data; D. Church and V. Schneider for helpful discussions about GRCh38; and many individuals for providing feedback on the current version and previous versions of our calls. Selected commercial equipment, instruments or materials are identified to specify the adequacy of experimental conditions or reported results. Such identification does not imply recommendation or endorsement by the National Institute of Standards, nor does it imply that the equipment, instruments or materials identified are necessarily the best available for the purpose. C.X. and S.S. were supported by the Intramural Research Program of the National Library of Medicine, National Institutes of Health. J.Z., M.S., N.O. and J.W. were supported by the National Institute of Standards and Technology and an interagency agreement with the Food and Drug Administration.

Author contributions

J.M.Z., L.T., N.D.O., J.W. and M.S. wrote the manuscript. J.M.Z., J.M., F.M.D., N.D.O., J.W., M.S. and H.P. designed and implemented the integration process. H.H., J.M. and J.M.Z. analyzed and integrated the 10x Genomics data. R.T., J.M. and J.M.Z. analyzed and integrated the Complete Genomics data. S.A.I., L.T., F.M.D., J.M. and J.M.Z. designed and implemented the phasing and robust trio analysis. C.Y.M., J.M. and J.M.Z. designed and implemented the robust GRCh38 liftover analysis. C.X. and S.S. managed and analyzed data. All authors contributed to GIAM discussions planning this work.

Competing interests

R.T. is an employee of, and holds stock in, Invitae. H.H. was an employee of 10x Genomics. S.A.I. and L.T. are employees of Real Time Genomics. C.Y.M. is an employee of Verily Life Sciences and Google.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41587-019-0074-6>.

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

Correspondence and requests for materials should be addressed to J.M.Z.

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

This is a US government work and not under copyright protection in the US; foreign copyright protection may apply, 2019

Methods

Sequencing datasets. In contrast to our previous integration process¹, which used sequencing data that were generated from DNA from various growths of the Coriell cell line GM12878, most of the sequencing data used in this work were generated from NIST RMs, which are DNA-extracted from a single large batch of cells that was mixed before aliquoting. These datasets, except for the HG001 SOLiD datasets, are described in detail in a GIAM data publication⁷.

For these genomes, we used datasets from several technologies and library preparations:

1. ~300× paired-end, whole-genome sequencing with 2 × 148 bp reads of 550 bp insert size from the HiSeq 2500 in Rapid Mode with v1 sequencing chemistry (HG001 and Ashkenazi Jewish (AJ) trio); the Chinese parents were similarly sequenced to ~100× coverage.
2. ~300× paired-end, whole-genome sequencing with 2 × 250 bp reads with ~50 bp insert size from the HiSeq 2500 in Rapid Mode with v2 sequencing chemistry (Chinese son).
3. ~45× paired-end, whole-genome sequencing with 2 × 250 bp reads with ~400 bp insert size from the HiSeq 2500 in Rapid Mode with v2 sequencing chemistry (AJ trio).
4. ~15× mate-pair, whole-genome sequencing with 2 × 100 bp reads with ~6,000 bp insert size from the HiSeq 2500 in high throughput mode with v2 sequencing chemistry (AJ trio and Chinese son).
5. ~100× paired-end, 2 × 29 bp whole-genome sequencing from Complete Genomics v2 chemistry (all genomes).
6. ~1,000× single-end exome sequencing from Ion PI Sequencing 200 Kit v.4 (HG001, AJ trio and Chinese son).
7. Whole-genome sequencing from SOLiD 5500 W. For HG001, two whole-genome sequencing datasets were generated: ~12 × 2 × 50 bp paired-end and ~12 × 75 bp single-end with error-correction chemistry. For the AJ son and Chinese son, ~42 × 75 bp single-end sequencing (without error-correction chemistry) was generated.
8. 10× Genomics Chromium whole-genome sequencing (~25× per haplotype from HG001 and AJ son, ~9× per haplotype from AJ parents and ~13× per haplotype from Chinese parents). These are the only data not derived from the NIST RM batch of DNA, because longer DNA from the cell lines resulted in better libraries.

Implementation of analyses and source code. Most analyses were performed using apps or applets on the DNAnexus data analysis platform, except for mapping of all datasets and variant calling for Complete Genomics and Ion exome, since these steps were performed previously. The apps and applets used in this work are included in GitHub (<https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2>). They run on an Ubuntu 12.04 machine on Amazon Web Services EC2. The apps and applets are structured thus:

1. `dxapp.json` specifies the input files and options, output files and any dependencies that can be installed via the Ubuntu command `apt`.
2. `src/code.sh` contains the commands that are run.
3. `resources/` contains compiled binary files, scripts and other files that are used in the applet.

The commands were run per chromosome in parallel using the DNAnexus command line interface, and these commands are listed at https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2/DNAnexusCommands/batch_processing_commands.

Illumina analyses. The Illumina fastq files were mapped using `novalign -d <reference.ndx> -f <read1.fastq.gz> <read2.fastq.gz> -F STDfQ --Q2off -t 400 -o SAM -c 10` (v.3.02.07 from Novocraft Technologies), and resulting BAM files were created, sorted and indexed with `samtools v.0.1.18`. Bam files are available under:

ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/NA12878/NIST_NA12878_HG001_HiSeq_300x/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_HiSeq_HG002_Homogeneity-10953946/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_Illumina_2x250bps/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG002_NA24385_son/NIST_Stanford_Illumina_6kb_matepair/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_HiSeq_HG003_Homogeneity-12389378/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_Illumina_2x250bps/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG003_NA24149_father/NIST_Stanford_Illumina_6kb_matepair/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_HiSeq_HG004_Homogeneity-14572558/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_Illumina_2x250bps/

ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/HG004_NA24143_mother/NIST_Stanford_Illumina_6kb_matepair/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG005_NA24631_son/HG005_NA24631_son_HiSeq_300x/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG005_NA24631_son/NIST_Stanford_Illumina_6kb_matepair/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG006_NA24694_huCA017E_father/NA24694_Father_HiSeq100x/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG006_NA24694_huCA017E_father/NIST_Stanford_Illumina_6kb_matepair/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG007_NA24695_hu38168_mother/NA24695_Mother_HiSeq100x/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/HG007_NA24695_hu38168_mother/NIST_Stanford_Illumina_6kb_matepair/

Variances were called using both GATK HaplotypeCaller v.3.5 (refs. ^{23,24}) and Freebayes 0.9.20 (ref. ²⁵) with high-sensitivity settings. Specifically, for HaplotypeCaller, special options were ‘-stand_call_conf 2 -stand_emit_conf 2 -A BaseQualityRankSumTest -A ClippingRankSumTest -A Coverage -A FisherStrand -A LowMQ -A RMSMappingQuality -A ReadPosRankSumTest -A StrandOddsRatio -A HomopolymerRun -A TandemRepeatAnnotator’. The gVCF output was converted to variant call format (VCF) using GATK Genotype gVCFs for each sample independently. For Freebayes, special options were ‘-F 0.05 -m 0 --genotype-qualities’.

For Freebayes calls, GATK CallableLoci v.3.5 was used to generate callable regions with at least 20 reads with mapping quality of at least 20 (to exclude regions where heterozygous variants may be missed), and with coverage less than twice the median (to exclude regions likely to be duplicated or to have mis-mapped reads). Because parallelization of Freebayes infrequently causes conflicting variant calls at the same position, two variants at the same position are removed from the callable regions.

For GATK calls, the gVCF output from GATK was used to define the callable regions. In general, reference regions and variant calls with GQ < 60 were excluded from the callable regions, excluding 50 bp on either side of low-GQ reference regions and low-GQ variant calls. Exclusion of 50 bp minimizes artifacts introduced when integrating partial complex variant calls. The exception to this rule is that reference assertions with GQ < 60 are ignored if they are within 10 bp of the start or end of an indel with GQ > 60, because following manual inspection, GATK often calls some reference bases with a low GQ near true indels, even when the reference bases should have high GQ, and exclusion of 50 bp regions around these bases excluded many true benchmark indels. The gVCF from GATK is used rather than CallableLoci, because it provides a sophisticated interrogation of homopolymers and tandem repeats and excludes regions if insufficient reads completely cross the repeats.

Complete Genomics analyses. Complete Genomics data was mapped and variants were called using v.2.5.0.33 of the standard Complete Genomics pipeline²⁶. Only the `vcfBeta` file was used in the integration process, because this contains both called variants and ‘no call’ regions similar to gVCF. `vcfBeta` files are available under:

ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/NA12878/analysis/CompleteGenomics_RMDNA_11272014/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/analysis/CompleteGenomics_RefMaterial_SmallVariants_CGAtools_08082014/
ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/ChineseTrio/analysis/CompleteGenomics_HanTrio_RMDNA_08082014/

A python script from Complete Genomics (`vcf2bed.py`) was used to generate callable regions, which exclude those with no calls or partial no calls in the `vcfBeta`. To minimize integration artifacts around partial complex variant calls, 50 bp padding was subtracted from both sides of callable regions using `bedtools slopBed`. In addition, the `vcfBeta` file was modified to remove unnecessary lines and to fill in the `FILTER` field using a python script from Complete Genomics (`vcfBeta_to_VCF_simple.py`). This process was performed on DNAnexus; an example command for chr20 is as follows:

```
dx run GLAB/Workflow/integration-prepare-cg -ivcf_in=/NA12878/Complete_Genomics/vcfBeta-GS000025639-ASM.vcf.bz2 -ichrom=20 --destination=/NA12878/Complete_Genomics/
```

Ion exome analyses. The Ion exome data BaseCalling and alignment were performed on a Torrent Suite v.4.2 server (ThermoFisher Scientific). Variant calling was performed using Torrent Variant Caller v.4.4, and the `TSVC_variants_defaultlowsetting.vcf` was used as a sensitive variant call file.

GATK CallableLoci v.3.5 was used to generate callable regions with at least 20 reads with mapping quality of at least 20.

These callable regions were intersected with the targeted regions `bed` file for the Ion Ampliseq exome assay available at ftp://ftp-trace.ncbi.nlm.nih.gov/trace/ftp/data/AshkenazimTrio/analysis/IonTorrent_TVC_03162015/AmpliseqExome.20141120_effective_regions.bed. In addition, 50 bp on either side of compound heterozygous sites were removed from the callable regions, and these sites were removed from the `vcf` to avoid artifacts around homopolymers. This process was performed on DNAnexus, with the command for chr20:

```
dx run GIAB/Workflow/integration-prepare-ion -ivcf_in=/NA12878/
Ion_Torrent/TSVC_variants_defaultlowsetting.vcf -icallableloci=/NA12878/
Ion_Torrent/callableLoci_output/HG001_20_hs37d5_IonExome_callableloci.
bed -itargetsbed=/NA12878/Ion_Torrent/AmpliseqExome.20141120_effective_
regions.bed -ichrom=20 --destination=/NA12878/Ion_Torrent/Integration_
prepare_ion_output/
```

SOLiD analyses. The SOLiD xsq files were mapped with LifeScope v.2.5.1 (ThermoFisher Scientific). Bam files are available under:

```
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NIST_NA12878_
HG001_SOLID5500W/
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/
```

Variants were called using GATK HaplotypeCaller v.3.5 with high-sensitivity settings. Specifically, for HaplotypeCaller, special options were ‘-stand_call_conf 2 -stand_emit_conf 2 -A BaseQualityRankSumTest -A ClippingRankSumTest -A Coverage -A FisherStrand -A LowMQ -A RMSMappingQuality -A ReadPosRankSumTest -A StrandOddsRatio -A HomopolymerRun -A TandemRepeatAnnotator’. The gVCF output was converted to VCF using GATK Genotype gVCFs for each sample independently.

For SOLiD, all regions were considered ‘not callable’ because biases are not sufficiently well understood, so SOLiD provided support only from an additional technology when finding training sites and when annotating the benchmark vcf.

10x Genomics analyses. The 10x Genomics Chromium fastq files were mapped and reads were phased using LongRanger v.2.0 (10x Genomics)²¹. As a new technology, variant calls are integrated in a conservative manner, requiring clear support for a homozygous variant or reference call from reads on each haplotype. Bamtools was used to split the bam file into two separate bam files with reads from each haplotype (HP tag values 1 and 2), ignoring reads that were not phased. GATK HaplotypeCaller v.3.5 was used to generate gvcf files from each haplotype separately. A custom perl script was used to parse the gvcf files, excluding regions with read depth (DP) < 6, DP > 2 × median coverage, heterozygous calls on a haplotype or homozygous reference or variant calls where the likelihood was < 99 for homozygous variant or reference, respectively (https://github.com/jzook/genome-data-integration/blob/master/NISTv3.3.2/DNAexusApplets/integration-prepare-10X-v3.3-anyref/resources/usr/bin/process_10X_gvcf.pl). The union of these regions from both haplotypes plus 50 bp on either side was excluded from the callable regions.

Exclusion of challenging regions for each callset. As an enhancement in our new integration methods (v.3.3+), we now exclude difficult regions per callset rather than at the end for all callsets. We used both previous knowledge and manual inspection of differences between callsets to determine which regions should be excluded from each callset. For example, we exclude tandem repeats and perfect or imperfect homopolymers > 10 bp for callsets with reads < 100 bp or from technologies that use PCR (all methods except Illumina PCR-free). We exclude segmental duplications and regions homologous to the decoy hs37d5 or ALT loci from all short-read methods (all methods except 10x Genomics-linked reads). The regions excluded from each callset are specified in the ‘CallsetTables’ input into the integration process, where a 1 in a column indicates that the bed file regions are excluded from that callset’s callable regions, and a 0 indicates that the bed file regions are not excluded for that callset. Callset tables are available at <https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2/CallsetTables>.

Integration process to form benchmark calls and regions. Our integration process is summarized in Fig. 1 and is detailed in the outline in Supplementary Note 10 and in diagrams in Supplementary Figs. 2, 3 and 4. Similar to our previous work for v.2.18/v.2.19 (ref. 1), the first step in our integration process is to use preliminary ‘concordant’ calls to train a machine learning model that finds outliers (Supplementary Fig. 5). Callsets with a genotype call that has outlier annotations are not trusted. For training variants, we use genotype calls that are supported by at least two technologies, not including sites if another dataset contradicts the call or is missing the call and that call is within that callset’s callable regions. To do this, we normalize variants using vcflib vcflib primitives and then generate a union vcf using vcflib vcfcombine. The union vcf contains separate columns from each callset, and we annotate the union vcf with vcflib vcfannotate to indicate whether each call falls outside the callable bed file from each dataset.

To find outliers, annotations from the vcf INFO and FORMAT fields, and which tail of the distribution that we expected to be associated with bias for each callset, were selected using expert knowledge (files describing annotations for each caller are at <https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2/AnnotationFiles>). Then, we used a simple one-class model that treats each annotation independently and finds the cutoff (for single tail) or cutoffs (for both tails) outside which ($5a^{-1}$)% of the training calls lie, where a is the number of annotations for each callset. For each callset, we find sites that fall outside the cutoffs or that are already filtered, and we generate a bed file that contains the call

with 50 bp added to either side to account for different representations of complex variants. We again annotate the union vcf with this ‘filter bed file’ from each callset, which we next use in addition to the ‘callable regions’ annotations (Fig. 1b).

To generate the benchmark calls, we run the same integration script with the new union vcf annotated with both callable regions from each dataset and ‘filtered regions’ from each callset (Fig. 1c). The integration script outputs benchmark calls that meet all of the following criteria (in the vcf file for each genome ending in _v3.3.2_all.vcf.gz, an appropriate output vcf FILTER status is given for sites that don’t meet each criterion, summarized in Supplementary Table 3):

1. Genotypes agree between all callsets for which the call is callable and not filtered. This includes ‘implied homozygous reference calls’, where a site is missing from a callset and it is within that callset’s callable regions and no filtered variants are within 50 bp (otherwise, FILTER status ‘discordantunfiltered’, for example, (2) in Fig. 1c).
2. The call from at least one callset is callable and not filtered (otherwise, FILTER status ‘allfilteredbutagree’ or ‘allfilteredanddisagree’, for example, (4) in Fig. 1c).
3. The sum of genotype qualities for all datasets supporting this genotype call is > 70. This sum includes only the first callset from each dataset and includes all datasets supporting the call, even if these are filtered or outside the callable bed (otherwise, FILTER status ‘GQlessthan70’).
4. If the call is a homozygous reference, then no filtered calls at the location can be indels > 9 bp (otherwise, FILTER status ‘questionableindel’). This criterion was added since implied homozygous reference calls (that is, there is no variant call and it is in the callset’s callable regions) are sometimes unreliable for larger indels, because callsets will sometimes miss the evidence for large indels and not call a variant.
5. The site is not called only by Complete Genomics and completely missing from other callsets, since these sites tended to be systematic errors in repetitive regions following manual curation (otherwise, FILTER status ‘cgonly’).
6. For sites where the benchmark call is heterozygous, none of the filtered calls are homozygous variant since these are sometimes genotyping errors (otherwise, FILTER status ‘discordanthet’).
7. Heterozygous calls where the net allele balance across all unfiltered datasets is < 0.2 or > 0.8 when summing support for reference and alternate alleles (otherwise, FILTER status ‘alleleimbalance’).

To calculate the benchmark regions bed file, the following steps are performed:

- (1) Find all regions that are covered by at least one dataset’s callable regions bed file.
- (2) Subtract 50 bp on either side of all sites that could not be determined with high confidence.

We chose to use a simple, one-class model in this work in place of the more sophisticated Gaussian mixture model (GMM) used in v.2.19 to make the integration process more robust and reproducible. The way in which we used the GMM in v.2.18 and v.2.19, GATK Variant Quality Score Recalibration was not always able to fit the data, requiring manually modified parameters, making the integration process less robust and reproducible. At times the GMM also appeared to overfit the data, filtering sites for unclear reasons. The GMM also was designed for GATK, ideally for vcf files containing many individuals, unlike our single-sample vcf files. Our one-class model used in v.3.3.2 was more easily adapted to filter on annotations in the FORMAT and INFO fields in vcf files produced by the variety of variant callers and technologies used for each sample. Future work could be directed towards developing more sophisticated models that learn which annotations are more useful, how to define the callable regions and how annotations may have different values for different types of variants and sizes.

The integration process described in this section is implemented in the applet nist-integration-v3.3.2-anyref (<https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2/DNAexusApplets/nist-integration-v3.3.2-anyref>). The one-class filtering model to find outliers is a custom perl script (nist-integration-v3.3.2-anyref/resources/usr/bin/VcfOneClassFiltering_v3.3.pl). The integration and arbitration processes to (1) find genotype calls supported by two technologies for the training set and (2) find benchmark genotype calls are both implemented in a perl script (nist-integration-v3.3.2-anyref/resources/usr/bin/VcfClassifyUsingFilters_v3.3.pl). Each line in the benchmark vcf file has FORMAT fields (Supplementary Table 4) and INFO fields (Supplementary Table 5) that give information about the support from each input callset for the final call. Some variants in the benchmark vcf fall outside the benchmark bed file, because inclusion of these variants when doing the variant comparison helps account for differences in representation of complex variants between the benchmark vcf and the vcf being tested.

Candidate structural variants to exclude from benchmark regions. For HG001, we had previously used all HG001 variants from dbVar, which conservatively excluded about 10% of the genome because dbVar contained all variants previously submitted for HG001. Since then, several SV callsets from different technologies have been generated for HG001, and we now use these callsets rather than dbVar. Specifically, we include union of all calls < 1 Mbp in size from:

1. The union of several PacBio SV calling methods, including filtered sites (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NA12878_PacBio_MtSi-nai/NA12878.sorted.vcf.gz)
2. The PASSing calls from MetaSV, which looks for support from multiple types of Illumina calling methods (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify_Manuscript/Supplementary_Information/metasv_trio_validation/NA12878_svs.vcf.gz)
3. Calls with support from multiple technologies from svclassify (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify_Manuscript/Supplementary_Information/Personalis_1000_Genomes_deduplicated_deletions.bed and ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify_Manuscript/Supplementary_Information/Spiral_Genetics_insertions.bed)

For the AJ trio, we use the union of calls from 11 callers using five technologies. Specifically, we include all calls >50 bp and <1 Mbp in size from:

1. PacBio callers: sniffles, Parliament, CSHL assembly, SMRT-SV dip, and Multi-breakSV
2. Illumina callers: Cortex, Spiral Genetics, Jitterbug, and Parliament
3. BioNano haplo-aware calls
4. 10x GemCode
5. Complete Genomics highConfidenceSvEventsBeta and mobileElementInsertionsBeta

For the Chinese trio, we use the union of calls from 11 callers using five technologies. Specifically, we include all calls >50 bp and <1 Mbp in size from:

1. Illumina callers: GATK-HC and Freebayes
2. Complete Genomics highConfidenceSvEventsBeta and mobileElementInsertionsBeta
3. For HG006 and HG007, we also exclude calls from MetaSV in any individual in the trio

For deletions we add 100 bp to either side of the called region, and for all other calls (mostly insertions) we add 1,000 bp to either side of the called region. This padding helps to account for imprecision in the called breakpoint, complex variation around the breakpoints and potential errors in large tandem duplications that are reported as insertions. For GRCh38, we remap the GRCh37 SV bed file to GRCh38 using NCBI remap.

GRCh38 integration. To develop benchmark calls for GRCh38, we use methods similar to those for GRCh37 except for data that could not be realigned to GRCh38. We are able to map reads and call variants directly on GRCh38 for Illumina and 10x, but native GRCh38 pipelines were not available for Complete Genomics, Ion exome and SOLiD data. For Illumina and 10x data, variant calls were made similarly to GRCh37 but from reads mapped to GRCh38 with decoy but no alternative loci (ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh38/seqs_for_alignment_pipelines.ucsc_ids/GCA_000001405.15_GRCh38_no_alt_plus_hs38d1_analysis_set.fna.gz). For Complete Genomics, Ion exome and SOLiD data, vcf and callable bed files were converted from GRCh37 to GRCh38 using the tool GenomeWarp (<https://github.com/verilylifesciences/genomewarp>). This tool converts vcf and callable bed files in a conservative and sophisticated manner, accounting for base changes that were made between the two references. Modeled centromere (genomic_regions_definitions_modeledcentromere.bed) and heterochromatin (genomic_regions_definitions_heterochrom.bed) regions are explicitly excluded from the benchmark bed (available under ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/NISTv3.3.2/GRCh38/supplementary/Files/).

For HG001 Illumina and 10x data, variants were called using GATK HaplotypeCaller v.3.5 similarly to all GRCh37 genomes but, for the other genomes' Illumina and 10x datasets, the Sentieon haplotyper v.201611.rc1 was used instead. This tool was designed to give the same results more efficiently than GATK HaplotypeCaller v.3.5 except that it does not downsample reads in high-coverage regions, so that the resulting variant calls are deterministic²⁷.

Trio Mendelian analysis and phasing. Because variants are called independently in the three individuals, complex variants are sometimes represented differently amongst these; the majority of apparent Mendelian violations found by naive methods are in fact discrepant representations of complex variants in different individuals. We used a new method based on rtg-tools vcfeval to harmonize representation of variants across the three individuals, similar to the methods described in a recent publication²⁸. We performed a Mendelian analysis and phasing for the AJ trio using rtg-tools v.3.7.1. First, we harmonize the representation of variants across the trio using an experimental mode of rtg-tools vcfeval, so that different representations of complex variants do not cause apparent Mendelian violations. After merging the three individuals' harmonized vcfs into a multi-sample vcf, we change missing genotypes for individuals to homozygous reference, since we later subset by benchmark regions (missing calls are implied to be homozygous reference in our benchmark regions). Then, we use rtg-tools Mendelian to phase variants in the son and find Mendelian violations. Finally, we subset the Mendelian violations by the intersection of the benchmark bed files from all three individuals. We also generate new benchmark bed files for

each individual that exclude 50 bp on each side of Mendelian violations that are not apparent de novo mutations. This analysis and the phase transfer below is performed in the DNAnexus applet *trio-harmonize-mendelian* (<https://github.com/jzook/genome-data-integration/tree/master/NISTv3.3.2/DNAnexusApplets/trio-harmonize-mendelian>).

Phase transfer. Our integration methods supply only local phasing where GATK HaplotypeCaller (or Sentieon haplotyper) is able to phase the variants. For HG001/NA12878 and AJ son (HG002), pedigree-based phasing information can supply global phasing of maternal and paternal haplotypes.

For HG001, Real Time Genomics (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/analysis/RTG_Illumina_Segregation_Phasing_05122016/) and Illumina PG (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/platinum_genomes/2016-1.0/) have generated phased variant calls by analyzing her 17-member pedigree. The PG calls are phased in the order paternal/maternal, and the RTG phasing is not ordered, so the phasing from PG was added first. To start, we archive existing local phasing information to the IGT and IPS fields. Then, we use the rtg-tools vcfeval software phase-transfer mode to take the phasing first from PG and add it to the genotype (GT) field of our HG001 benchmark vcf. We then use the rtg-tools vcfeval's phase-transfer mode to take the phasing from RTG and add it to variants that were not phased by PG, flipping the phasing to the order paternal/maternal where necessary. Variants with phasing from PG or RTG are given a phase set (PS) = PATMAT. Variants that were not phased by either PG or RTG but are homozygous variant are given GT = 1|1 and PS = HOMVAR. Variants that were not phased by either PG or RTG but were phased locally by GATK HaplotypeCaller or Sentieon haplotyper are given the phased GT and PS from the variant caller. This phase transfer is performed in two steps with the applet *phase-transfer*.

For HG002, we apply the phasing from the trio analysis using rtg-tools vcfeval phase-transfer as above. Variants with trio-based phasing are given a PS = PATMAT. Variants that were not phased by the trio but are homozygous variant are given GT = 1|1 and PS = HOMVAR. Variants that were not phased by the trio but were phased locally by GATK HaplotypeCaller or Sentieon haplotyper are given the phased GT and PS from the variant caller. This phase transfer is performed in the applet *trio-harmonize-mendelian*.

Comparisons to other callsets. The Global Alliance for Genomics and Health formed a benchmarking team (<https://github.com/ga4gh/benchmarking-tools>) to standardize performance metrics and develop tools to compare different representations of complex variants. We have used one of these tools, vcfeval (<https://github.com/RealTimeGenomics/rtg-tools>) from rtg-tools-3.6.2 with --ref-overlap, to compare our benchmark calls to other vcfs. After performing the comparison, we subset the true-positives, false-positives and false-negatives by our benchmark bed file and then by the bed file accompanying the other vcf (if it has one). We then manually inspect alignments from a subset of the putative false-positives and -negatives and record whether we determine that our benchmark call is probably correct, if we understand why the other callset is incorrect, if the evidence is unclear, if it is in a homopolymer, and other notes.

For HG001, we compare to these callsets:

1. NISTv.2.18 benchmark calls and bed file that we published previously
2. PG 2016-1.0, with a modified benchmark bed file that excludes an additional 50 bp around uncertain variants or regions. This padding eliminates many locations where PG calls only part of a complex or compound heterozygous variant

In addition, we use the hap.py + vcfeval benchmarking tool (<https://github.com/Illumina/hap.py>) developed by the GA4GH Benchmarking Team and implemented on the precisionFDA website (<https://precision.fda.gov/apps/app-F187Zbj0qXjB85Yq2B6P61zb>) for use in the precisionFDA 'Truth Challenge'. We modified the tool used in the challenge to stratify performance by additional bed files available at <https://github.com/ga4gh/benchmarking-tools/tree/master/resources/stratification-bed-files>, including bed files of 'easier' regions and those encompassing complex and compound heterozygous variants. The results of these comparisons, as well as the pipelines used to generate the calls, are shared in a Note on precisionFDA (<https://precision.fda.gov/notes/300-giab-example-comparisons-v3-3-2>). These precisionFDA results can be accessed immediately by requesting a free account on precisionFDA.

The callsets with lowest false-positive and -negative rates for SNVs or indels from the precisionFDA Truth Challenge were compared to the v.3.3.2 benchmark calls for HG001. From each comparison result, ten putative false-positives or -negatives were randomly selected for manual inspection to assess whether they were, in fact, errors in each callset.

Integration with only Illumina and 10x Genomics WGS. To assess the impact of using fewer datasets on the resulting benchmark vcf and bed files, we performed integration for chromosome 1 in GRCh37 using only Illumina 300x 2 × 150 bp WGS and 10x Genomics Chromium data for HG001. We compared these calls to PG as described above for v.3.3.2 calls, and we manually inspected all differences with PG that were not in v.3.3.2.

Differences between old and new integration methods. The new integration methods differ from the previous GIAB calls (v.2.18 and v.2.19) in several ways, both in the data used and the integration process and heuristics:

(1) Only newer datasets were used, which were generated from the NIST RM 8398 batch of DNA (except for 10x Genomics, which used longer DNA from cells).

(2) Mapping and variant calling algorithms designed specifically for each technology were used to generate sensitive variant callsets where possible: novoalign + GATK-haplotypecaller and Freebayes for Illumina, the vcfBeta file from the standard Complete Genomics pipeline, tmap + TVC for Ion exome and LifeScope + GATK-HC for SOLiD. This is intended to minimize bias towards any particular bioinformatics toolchain.

(3) Rather than forcing GATK to call genotypes at candidate variants in the bam files from each technology, we generate sensitive variant callsets and a bed file that describes the regions that were callable from each dataset. For Illumina GATK, we used the GATK-HC gVCF output to find regions with GQ > 60. For Illumina Freebayes, we used GATK callable loci to find regions with at least 20 reads with MQ ≥ 20 and with coverage less than twice the median. For Complete Genomics, we used the callable regions defined by the vcfBeta file and excluded ±50 bp around any no-called or half-called variant. For Ion, we intersected the exome targeted regions with the output of GATK CallableLoci for the bam file (requiring at least 20 reads with MQ ≥ 20). Due to the shorter reads and low coverage for SOLiD, this was used only to confirm variants, so no regions were considered callable.

(4) A new file with putative structural variants was used to exclude potential errors around SVs. For HG001, these were SVs derived from multiple PacBio callers (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NA12878_PacBio_MtSinai/) and multiple integrated Illumina callers using MetaSV (ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify_Manuscript/Supplementary_Information/metasy_trio_validation/). These comprise a notably lower fraction of the calls and genome (~4.5%) than the previous bed, which was a union of all dbVar calls for HG001 (~10%). For the AJ trio, the union of fewer than ten submitted SV callsets from Illumina, PacBio, BioNano and Complete Genomics from all members of the trio combined was used to exclude potential SV regions. For the Chinese trio, only Complete Genomics and GATK-HC and Freebayes calls >49 bp and surrounding regions were excluded due to the lack of available SV callsets for this genome at this time, which may have resulted in a higher error rate in this genome. The SV bed files for each genome are included in the Supplementary Files directory.

(5) To eliminate some errors from v.2.18, homopolymers >10 bp in length, including those interrupted by one nucleotide different from the homopolymer, are excluded from all input callsets except PCR-free GATK HaplotypeCaller callsets; for these callsets, we include only sites with confident genotype calls where HaplotypeCaller has ensured that sufficient reads entirely encompass the repeat.

(6) A bug that caused nearby variants to be missed in v.2.19 is fixed in the new calls.

(7) The new vcf contains variants outside the benchmark bed file. This enables more robust comparison of complex variants or nearby variants that are near the boundary of the bed file. It also allows the user to evaluate concordance outside the benchmark regions, but these concordance metrics should be interpreted with great care.

(8) We now supply global phasing information from pedigree-based calls for HG001, trio-based phasing for the AJ son and local phasing information from GATK-HC for the other genomes.

(9) We use phased reads to make variant calls from 10x Genomics, conservatively requiring at least six reads from both haplotypes, coverage less than twice the median on each haplotype and clear support for either the reference allele or variant allele in each haplotype.

Statistical methods for performance metrics. The Wilson method in the R Hmisc binconf function was used to calculate 95% binomial confidence intervals for the recall and precision statistics in Supplementary Table 2.

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

Data availability

Raw sequence data were previously published in *Scientific Data* (<https://doi.org/10.1038/sdata.2016.25>) and were deposited in the NCBI SRA with the accession codes SRX1049768–SRX1049855, SRX847862–SRX848317, SRX1388368–SRX1388459, SRX1388732–SRX1388743, SRX852932–SRX852936, SRX847094, SRX848742–SRX848744, SRX326642, SRX1497273 and SRX1497276. 10x Genomics Chromium bam files used are available at ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/10XGenomics_ChromiumGenome_LongRanger2.0_06202016/. The benchmark vcf and bed files resulting from work in this manuscript are available in the NISTv3.3.2 directory under each genome on the GIAB FTP release folder <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/> and, in the future, updated calls will be in the ‘recent’ directory under each genome. The data used in this manuscript and other datasets for these genomes are available at <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/>, as well as in NCBI BioProject No. PRJNA200694.

Code availability

All code for analyzing genome sequencing data to generate benchmark variants and regions developed for this manuscript is available in a GitHub repository at <https://github.com/jzook/genome-data-integration>. Publicly available software used to generate input callsets includes novoalign v.3.02.07, samtools v.0.1.18, GATK v.3.5, Freebayes v.0.9.20, Complete Genomics tools v.2.5.0.33, Torrent Variant Caller v.4.4, LifeScope v.2.5.1, LongRanger v.2.0, GenomeWarp, rtg-tools v.3.7.1 and Sentieon v.201611.rc1.

References

23. DePristo, M. A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
24. McKenna, A. et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
25. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. Preprint at <https://arxiv.org/abs/1207.3907v2> (2012).
26. Drmanac, R. et al. Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays. *Science* **327**, 78–81 (2010).
27. Kendig, K. et al. Computational performance and accuracy of Sentieon DNaseq variant calling workflow. Preprint at *bioRxiv* 396325 <https://doi.org/10.1101/396325> (2018).
28. Toptaş, B. Ç., Rakocevic, G., Kómár, P. & Kural, D. Comparing complex variants in family trios. *Bioinformatics* <https://doi.org/10.1093/bioinformatics/bty443> (2018).

Reporting Summary

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

Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

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

Our web collection on [statistics for biologists](#) may be useful.

Software and code

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

Data collection

No software used for data collection.

Data analysis

All code for analyzing genome sequencing data to generate high-confidence variants and regions developed for this manuscript are available in a GitHub repository at <https://github.com/jzook/genome-data-integration>. Publicly available software used to generate input callsets includes novoalign version 3.02.07, samtools version 0.1.18, GATK v3.5, Freebayes 0.9.20, Complete Genomics tools v2.5.0.33, Torrent Variant Caller v4.4, LifeScope v2.5.1, LongRanger v2.1, GenomeWarp, rtg-tools v3.7.1, and sentieon version 201611.rc1.

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

Data

Policy information about [availability of data](#)

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

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

Raw sequence data were previously published in Scientific Data (DOI: 10.1038/sdata.2016.25), and were deposited in the NCBI SRA with the accession codes SRX1049768 to SRX1049855, SRX847862 to SRX848317, SRX1388368 to SRX1388459, SRX1388732 to SRX1388743, SRX852932 to SRX852936, SRX847094, SRX848742 to SRX848744, SRX326642, SRX1497273, and SRX1497276. 10x Genomics Chromium bam files used are at ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/10XGenomics_ChromiumGenome_LongRanger2.0_06202016/. The high-confidence vcf and bed files resulting from work in this manuscript are available in the NISTv3.3.2 directory under each genome on the GIAB FTP release folder <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/>, and in the future updated calls will be in the “recent” directory under each genome. The data used in this manuscript and other datasets for these genomes are available in <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/>, and in the NCBI BioProject PRJNA200694.

Field-specific reporting

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

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/authors/policies/ReportingSummary-flat.pdf](https://www.nature.com/authors/policies/ReportingSummary-flat.pdf)

Life sciences study design

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

Sample size	7 human genomes were characterized as benchmarks because these were all of the samples currently chosen for extensive characterization by the Genome in a Bottle Consortium
Data exclusions	No data excluded
Replication	Overall statistics of the benchmark sets were compared across all seven GIAB genomes to determine reproducibility across samples.
Randomization	Randomization is not relevant to our study, as there were not distinct experimental groups
Blinding	These benchmark samples are an open science resource, so no information is blinded.

Reporting for specific materials, systems and methods

Materials & experimental systems

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

Methods

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

Eukaryotic cell lines

Policy information about [cell lines](#)

Cell line source(s)	Coriell NIGMS Cell Line Repository (GM24385, GM24149, GM24143, GM24631, GM24694, GM24695, GM12878)
Authentication	Whole genome sequencing and variant calling was performed on all specimens
Mycoplasma contamination	All cell lines tested negative for mycoplasma contamination

Commonly misidentified lines
(See [ICLAC](#) register)

No commonly misidentified cell lines were used.