

A framework for variation discovery and genotyping using next-generation DNA sequencing data

Mark A DePristo¹, Eric Banks¹, Ryan Poplin¹, Kiran V Garimella¹, Jared R Maguire¹, Christopher Hartl¹, Anthony A Philippakis¹⁻³, Guillermo del Angel¹, Manuel A Rivas^{1,4}, Matt Hanna¹, Aaron McKenna¹, Tim J Fennell¹, Andrew M Kernytsky¹, Andrey Y Sivachenko¹, Kristian Cibulskis¹, Stacey B Gabriel¹, David Altshuler^{1,3,4} & Mark J Daly^{1,3,4}

Recent advances in sequencing technology make it possible to comprehensively catalog genetic variation in population samples, creating a foundation for understanding human disease, ancestry and evolution. The amounts of raw data produced are prodigious, and many computational steps are required to translate this output into high-quality variant calls. We present a unified analytic framework to discover and genotype variation among multiple samples simultaneously that achieves sensitive and specific results across five sequencing technologies and three distinct, canonical experimental designs. Our process includes (i) initial read mapping; (ii) local realignment around indels; (iii) base quality score recalibration; (iv) SNP discovery and genotyping to find all potential variants; and (v) machine learning to separate true segregating variation from machine artifacts common to next-generation sequencing technologies. We here discuss the application of these tools, instantiated in the Genome Analysis Toolkit, to deep whole-genome, whole-exome capture and multi-sample low-pass (~4×) 1000 Genomes Project datasets.

Recent advances in next-generation sequencing (NGS) technology now provide the first cost-effective approach to large-scale resequencing of human samples for medical and population genetics. Projects such as the 1000 Genomes Project¹ (1KG), The Cancer Genome Atlas and numerous large medically focused exome sequencing projects² are underway in an attempt to elucidate the full spectrum of human genetic diversity¹ and the complete genetic architecture of human disease. The ability to examine the entire genome in an unbiased way will make possible comprehensive searches for standing variation in common disease and mutations underlying linkages in Mendelian disease³, as well as spontaneously arising variation for which no gene-mapping shortcuts are available (for example, somatic mutations in cancer⁴⁻⁶ and *de novo* mutations⁷ (Conrad, D.F. *et al.* unpublished data) in autism and schizophrenia).

Many capabilities are required to obtain a complete and accurate record of the variation from NGS from sequencing data. Mapping

reads to the reference genome⁸⁻¹¹ is a first critical computational challenge whose cost necessitates that each read be aligned independently, guaranteeing that many reads spanning indels will be misaligned. The per-base quality scores, which convey the probability that the called base in the read is the true sequenced base¹², are quite inaccurate and co-vary with features like sequencing technology, machine cycle and sequence context¹³⁻¹⁵. These misaligned reads and inaccurate quality scores propagate into SNP discovery and genotyping, a general problem that becomes acute in projects with multiple sequencing technologies generated by many centers using rapidly evolving experimental processing pipelines, such as the 1000 Genomes Project.

Given well-mapped, aligned and calibrated reads, resolving even simple SNPs, let alone more complex variation such as multi-nucleotide substitutions, insertions and deletions, inversions, rearrangements and copy number variation, requires sensitive and specific statistical models^{8-11,15-25}. Separating true variation from machine artifacts as a result of the high rate and context-specific nature of sequencing errors is the outstanding challenge in NGS analysis. Previous approaches have relied on filtering SNP calls that have characteristics outside of their normal ranges, such as those occurring at sites with too much coverage^{17,19}, or by requiring non-reference bases to occur on at least three reads in both synthesis orientations²⁰. Though effective, such hard filters are frustratingly difficult to develop, require parameterization for each new dataset and are necessarily either restrictive (high specificity, as in the 1000 Genomes Project) or tolerant (high sensitivity, used in Mendelian disease studies, with concomitantly more false positives). Moreover, all of these challenges must be addressed within the context of a proliferation of sequencing technology platforms and study designs (for example, whole-genome shotgun, exome capture sequencing and multiple samples sequenced at shallow coverage), a point not tackled in previous work.

Here we present a single framework and the associated tools capable of discovering high-quality variation and genotyping individual samples using diverse sequencing machines and experimental designs (Fig. 1).

¹Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, Massachusetts, USA. ²Brigham and Women's Hospital, Boston, Massachusetts, USA. ³Harvard Medical School, Boston, Massachusetts, USA. ⁴Center for Human Genetic Research, Massachusetts General Hospital, Richard B. Simches Research Center, Boston, Massachusetts, USA. Correspondence should be addressed to M.A.D. (depristo@broadinstitute.org).

Received 27 August 2010; accepted 17 March 2011; published online 10 April 2011; doi:10.1038/ng.806

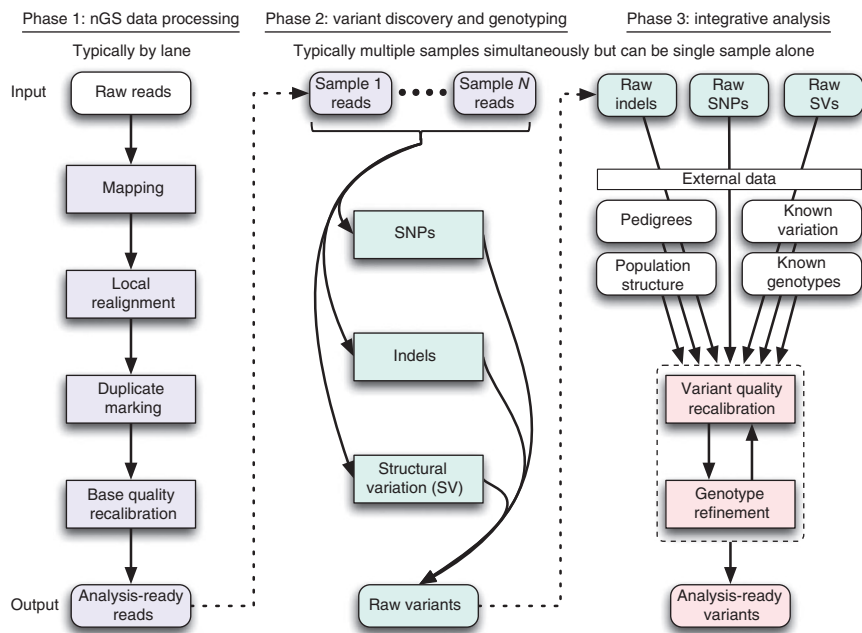


Figure 1 Framework for variation discovery and genotyping from next-generation DNA sequencing. See text for a detailed description.

We present several new methods addressing the challenges listed above in local realignment, base quality recalibration, multi-sample SNP calling and adaptive error modeling, which we apply to three prototypical NGS datasets (Table 1). In each dataset, we included CEPH individual NA12878 to show the consistency of results for this individual across all three datasets.

RESULTS

Below we describe a three-part conceptual framework (Fig. 1).

- Phase 1: raw read data with platform-dependent biases were transformed into a single, generic representation with well-calibrated base error estimates, mapped to their correct genomic origin and aligned consistently with respect to one another. Mapping algorithms placed reads with an initial alignment on the reference genome, either generated in, or converted to, the technology-independent SAM reference file format²⁴. Next, molecular duplicates were eliminated (Supplementary Note), initial alignments were refined by local realignment and then an empirically accurate per-base error model was determined.

- Phase 2: the analysis-ready SAM/BAM files were analyzed to discover all sites with statistical evidence for an alternate allele present among the samples including SNPs, short indels and copy number variations (CNVs). CNV discovery and genotyping methods, though part of this conceptual framework, are described elsewhere²⁵.

- Phase 3: technical covariates, known sites of variation, genotypes for individuals, linkage disequilibrium (LD), and family and population structure were integrated with the raw variant calls from phase 2 to separate true polymorphic sites from machine artifacts, and at these sites, high-quality genotypes were determined for all samples.

All components after initial mapping and duplicate marking were instantiated in the Genome Analysis Toolkit (GATK)²⁶.

Applying the analysis pipeline to HiSeq

Of the 2.83 billion non-N bases in the autosomal regions and chromosome X of the human reference genome, 2.72 billion bases (~96%) had sufficient coverage to call variants in the 101-bp paired-ended HiSeq data (Table 1). Even though the HiSeq reads were aligned with the gap-enabled BWA¹⁰, more than 15% of the reads that span known homozygous indels in NA12878 were misaligned (Supplementary Table 1). Realignment corrected 6.6 million of 2.4 billion total reads in 950,000 regions covering 21 Mb in the HiSeq data, eliminating 1.8 million loci with substantial accumulation of mismatching bases (Supplementary Table 2). The initial data-processing steps (phase 1) eliminated ~300,000 SNP calls, which is more than one fifth of the raw new calls, with quality metrics consistent with more than 90% of these SNPs being false positives (Table 2).

The initial 4.2 million confidently called non-reference sites included 99.7% and 99.5% of the HapMap3 and 1KG Trio sites, respectively, genotyped as non-reference in NA12878; at these variant sites, the sequencing and genotyping calls were concordant 99.9% of the time (Table 2). Variant quality score recalibration of these initial calls identified a tranche of SNPs with estimated false discovery rate (FDR) of <1%, containing 3.2 million known variants and 362,000 new variants, a 90% dbSNP rate, and transition/transversion (Ti/Tv) ratios of 2.15 and 2.05, respectively, consistent with our genome-wide expectations (Online Methods). Although the variant recalibrator removed ~595,000 total variants with a Ti/Tv ratio of ~1.2, it retained 99% of the HapMap3 and 97.3% of the 1KG Trio non-reference sites. The discordant sites have 100 times higher genotype discrepancy rates, suggesting that the sites themselves may be problematic. Almost all of the variants in the 1% tranche are already present in the even higher stringency 0.1% FDR tranche, and analysis of the 10% FDR tranche suggests that some more variants could be obtained, but at the cost of many more false positives.

Table 1 Next-generation DNA sequencing datasets analyzed

	HiSeq	Exome	Low-pass
Samples	NA12878	NA12878	NA12878 + 60 unrelated CEPH individuals
Sequencing technologies	Whole genome shotgun; Illumina HiSequation (2000) ¹⁷	Agilent exome hybrid capture ^{31,32} ; Illumina GenomeAnalyzer ¹⁷	Whole genome shotgun; Illumina GenomeAnalyzer ¹⁷ ; Life/SOLiD ³³ ; Roche/454 (ref. 19)
Coverage per sample	~60x	~150x; 93% of bases at >20x coverage	~4x
Read architecture	101 bp paired end	76/101 bp paired end	25, 36, 51, 76, ~250 (454) bp single and paired ends
Targeted area	2.85 Gb of autosomes and chr. X	28 Mb	2.85 Gb of autosomes and chr. X
Data set source	New, generated for this article	New, generated for this article	1000 Genomes Project
Aligner(s)	BWA ¹⁰	MAQ ⁹	MAQ ¹⁰ ; Corona Lite; SSAHA ¹²

Chr., chromosome.



Table 2 Raw to recalibrated, imputed SNP calls HiSeq, Exome and 61 sample low-pass datasets

Call set	Site discovery						Comparison to NA12878 variants			
	No. of SNPs				Ti/Tv		HM3 concordance		1KG concordance	
	All	Known	Novel	dbSNP (%)	Known	Novel	NR sensitivity	NRD rate	NR sensitivity	NRD rate
HiSeq										
Raw reads, all calls	4.43M	3.49M	941K	78.77	2.05	1.29	99.74	0.10	99.57	0.20
Unique to raw read calls	263K	37K	226K	13.95	1.37	0.70	0.02	37.97	0.09	12.64
Unique to +recal/+MSA calls	9.8K	1.8K	8.0K	18.08	1.38	1.39	0.00	18.18	0.00	9.93
+recal/+MSA, all calls	4.18M	3.45M	722K	82.71	2.06	1.57	99.72	0.09	99.48	0.19
Filtered by variant recalibration	595K	235K	360K	39.44	1.19	1.21	0.67	3.00	2.2	4.31
Final call set	3.58M	3.22M	362K	89.89	2.15	2.05	99.05	0.07	97.28	0.10
Low pass										
Raw reads, all calls	13.4M	6.5M	6.9M	48.77	2.05	1.13	83.97	20.34	80.45	22.53
Unique to raw read calls	670K	32K	638K	4.74	1.19	0.67	0.01	49.21	0.02	52.57
Unique to +recal/+MSA calls	45K	2.5K	42K	5.62	0.94	0.68	0.00	N/A	0.00	38.89
+recal/+MSA, all calls	12.8M	6.5M	6.3M	50.92	2.06	1.18	83.97	20.33	80.43	22.52
Filtered by variant recalibration	5.5M	706K	4.8M	12.84	1.31	1.01	0.95	26.54	3.44	32.91
Variant recalibrated call set	7.3M	5.8M	1.5M	79.7	2.18	2.05	Itemized below			
Sample variant calls for NA12878 only										
Variant recalibrated NGS reads only	2.44M	2.30M	140K	94.28	2.15	2.06	83.02	20.26	76.99	22.01
Recalibrated with Beagle imputation	3.20M	3.01M	191K	94.03	2.18	2.09	96.72	3.32	91.21	3.35
Exome capture										
Raw reads, all calls	18.9K	16.8K	2.1K	88.83	3.20	1.16	99.10	0.09	99.12	0.12
Unique to raw read calls	483	39	444	8.07	2.55	0.31	0.04	25.00	0.03	33.33
Unique to +recal/+MSA calls	81	40	41	49.38	3.44	1.73	0.01	0.00	0.04	16.67
+recal/+MSA, all calls	18.5K	16.8K	1.7K	90.77	3.20	1.61	99.07	0.08	99.13	0.11
Filtered by variant recalibration	1,274	609	665	47.8	1.85	0.84	0.59	N/A	0.76	N/A
Final call set	17.2K	16.2K	1,039	93.96	3.27	2.57	98.49	0.08	98.38	0.11

Part one of each section summarizes the impact of local realignment and base quality recalibration by comparing SNP calls on reads with raw quality scores and alignments to those made on the realigned, recalibrated reads. M, million; K, thousand.

Applying the analysis pipeline to 28-Mb exome capture

The raw data processing tools here eliminated ~450 new call sites from the raw call set, representing more than 20% of all the new calls, with a Ti/Tv of 0.30—fully consistent with all being false positives—and adding several sites present in HapMap3 and the 1KG Trio. The raw whole-exome data-call set, at ~150× coverage (Table 1), includes >99% of both the HapMap3 and 1KG Trio non-reference sites within the 28-Mb exome target region, with >99.8% genotype concordance at these sites. As with the HiSeq data, even with recalibration and local realignment, the Ti/Tv ratio of the new sites in the initial SNP calls indicates that more than 50% of these calls are false positives. Variant quality score recalibration, using only ~5,400 SNPs for training, identified a high-quality subset of calls that captured >98% of the HapMap3 and 1KG Trio sites in the target regions. The value of the tranches was more pronounced in the whole exome (Fig. 4d), where 900 of the 1,039 new calls come from tranches with FDRs under 1%, despite needing to reach into the 10% FDR tranche to include most true positive SNPs.

The HiSeq whole genome shotgun (WGS) and exome capture datasets differed drastically in their sequencing protocols (WGS versus hybrid capture), the sequencing machines (HiSeq versus Genome Analyzer) and the initial alignment tools (BWA¹⁰ versus MAQ⁹). Nevertheless, the exome call set is remarkably consistent the subset of calls from HiSeq that overlap the target regions of the hybrid capture protocol. Ninety-four percent of the HiSeq calls were also called in the final exome set sliced at 10% FDR (data not shown), and at these sites, the non-reference discrepancy rate was extremely low (<0.4%). Mapping differences between the aligners used for HiSeq (BWA) and exome (MAQ) datasets accounted for vast the majority of these discordant calls, with the remainder of the differences being because of limited coverage in the exome and only a small minority of

sites being because of differential SNP calling or variant quality score recalibration. Overall, despite the technical differences in the capture and sequencing protocols of the HiSeq and exome datasets, the data processing pipeline presented here uncovered a remarkably consistent set of SNPs in exomes with excellent genotyping accuracy.

Applying the analysis pipeline to low-pass (4×) sequencing

Multi-sample low-pass resequencing poses a major challenge for variant discovery and genotyping because there is so little evidence at any particular locus in the genome for any given sample (Table 1). Consequently, it is in precisely this situation, where there is little signal from true SNPs, that our data processing tools are most valuable, as can be seen from the progression of call sets in Table 2. Local realignment and base quality recalibration eliminated ~650,000 false-positive SNPs among 13 million sites, 4 times more sites than in the HiSeq dataset, with an aggregate Ti/Tv of 0.7. The initial low-pass CEU set includes over 13 million called sites among all individuals, of which nearly 7 million are new. NA12878 herself has 2.9 million variants, of which 430,000 are new. The 4× average coverage limits the sensitivity and concordance of this call set, with only 84% and 80% of HapMap3 and 1KG Trio sites, respectively, assigned a non-reference genotype in the NA12878 sample, both with a ~20% non-reference discrepancy (NRD) rate.

The variant quality recalibrator identified from the 13 million potential variants ~6 million known and 1.5 million new sites in tranches with 0.1% to 10% FDR. Figure 5a highlights several key features of the data: the allele frequency distribution of these calls closely matched the population genetics expectation, and the vast majority of HapMap3 and 1000 Genomes Project official CEU call sites were recovered, with the proportion nearing 100% for more common variant sites (Fig. 5a). Although we selected a 0.1% FDR

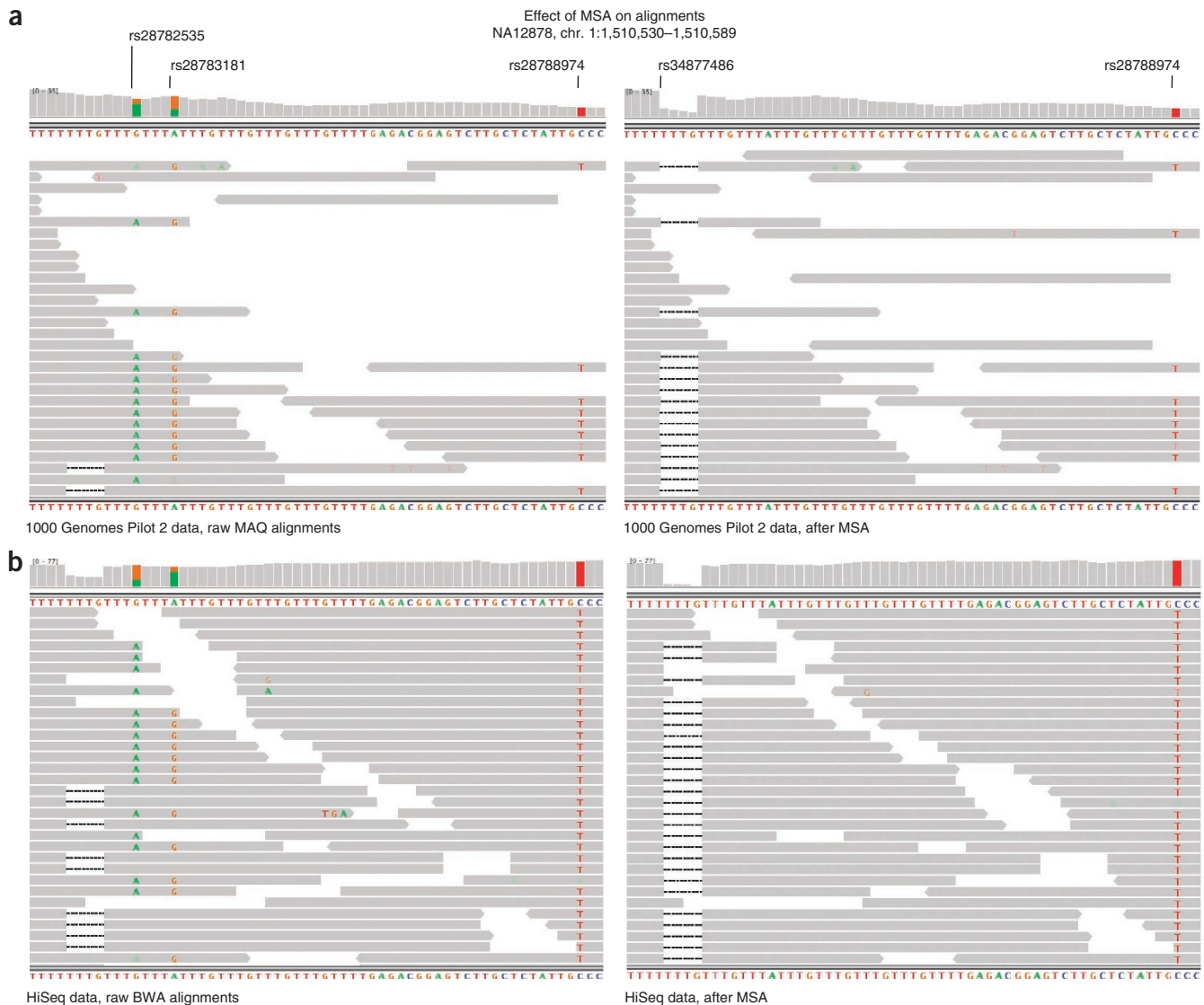


Figure 2 Integrative genomics viewer (IGV) visualization of alignments in region chr. 1: 1,510,530–1,510,589 from the Trio NA12878 Illumina reads from the 1000 Genomes Project (a) and NA12878 HiSeq reads before (left) and after (right) multiple sequence realignment (b). Reads are depicted as arrows oriented by increasing machine cycle; highlighted bases indicate mismatches to the reference: green, A; orange, G; red, T dashes, deleted bases a coverage histogram per base is shown above the reads. Both the 4-bp indel (rs34877486) and the C/T polymorphism (rs2878874) are present in dbSNP, as are the artifactual A/G polymorphisms (rs28782535 and rs28783181) resulting from the mis-modeled indel, indicating that these sites are common misalignment errors.

tranche for analysis here, which contains the bulk of HapMap3, 1KG Trio and HiSeq sites, there are another ~700,000 true sites that can be found in the 1% and 10% FDR tranches, albeit among many more false positives. This highest-quality tranche includes nearly all variants observed more than five times in the samples and 1.4 million new variants, with the SNPs in the tranches at 1% and 10% FDR generally occupying the lower alternate allele frequency range (Fig. 5b). The overall picture is clear: calling multiple samples simultaneously, even with only a handful of reads spanning a SNP for any given sample, enables one to detect the vast majority of common variant sites present in the cohort with a high degree of sensitivity.

Although the bulk properties of the 61-sample call set were good, we expected the low-pass 4x design to limit variation discovery and genotyping in each sample relative to deep resequencing. In the 61-sample call set, we discovered ~80% of the non-reference

sites in NA12878 according to the HapMap3, 1KG Trio and HiSeq call sets (Table 2). The ~20% of the missed variant sites from these three datasets had little to no coverage in the NA12878 sample in the low-pass data and, therefore, could not be assigned a genotype using only the NGS data, a general limitation of the low-pass sequencing strategy (Table 2 and Fig. 5c,d). The multi-sample discovery design, however, affords us the opportunity to apply imputation to refine and recover genotypes at sites with little or no sequencing data. Applying genotype-likelihood-based imputation with Beagle²⁷ to the 61-sample call set recovered an additional 15–20% of the non-reference sites in NA12878 that had insufficient coverage in the sequencing data (Table 2) as well as vastly improving genotyping accuracy (Fig. 5c,d).

We further characterized the quality of our low-pass call set as a function of the number of samples included during the discovery



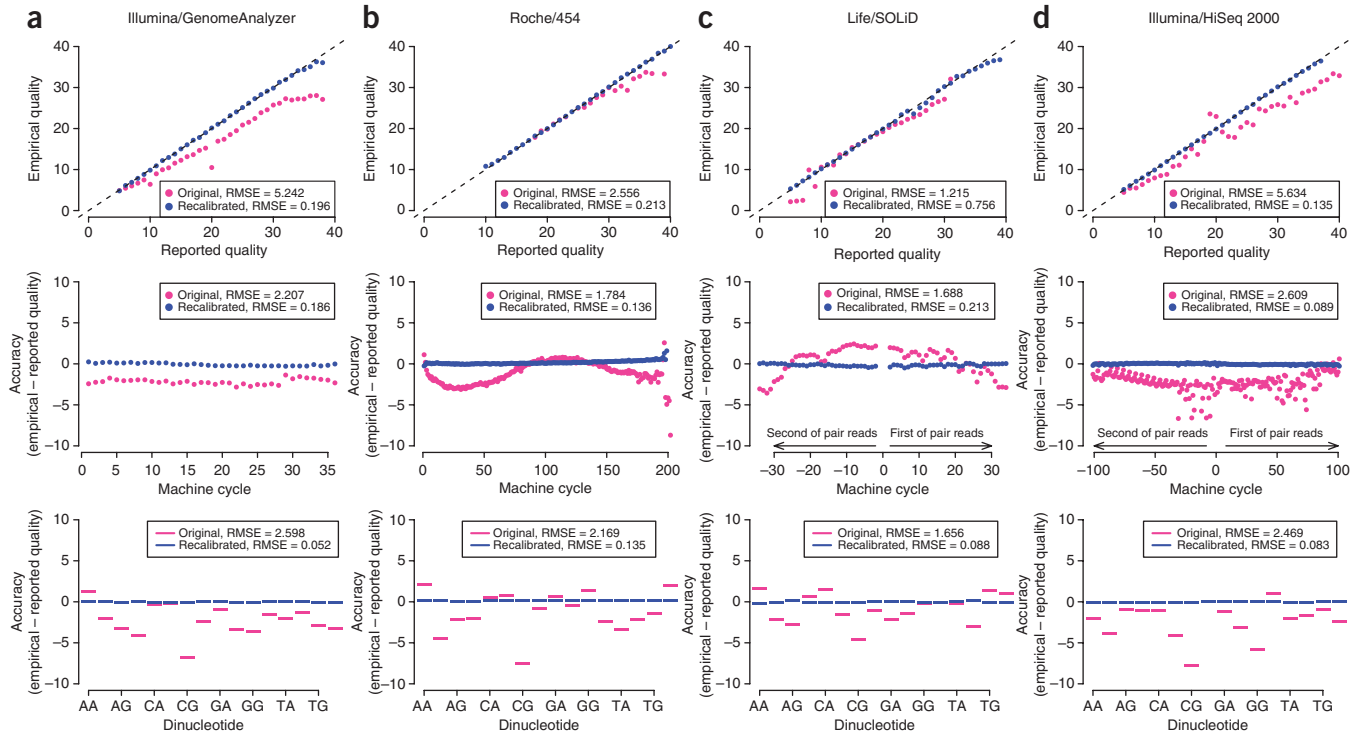


Figure 3 Raw (pink) and recalibrated (blue) base quality scores for NGS paired-end read sets of NA12878 of Illumina/GA (a), Roche/454 (b) and Life/SOLiD (c) lanes from the 1000 Genomes Project and Illumina/HiSeq (d). For each technology, the top panel shows reported base quality scores compared to the empirical estimates (Online Methods); the middle panel shows the difference between the average reported and empirical quality score for each machine cycle, with positive and negative cycle values given for the first and second read in the pair, respectively; and the bottom panel shows the difference between reported and empirical quality scores for each of the 16 genomic dinucleotide contexts. For example, the AG context occurs at all sites in a read where G is the current nucleotide and A is the preceding one in the read. Root-mean-square errors (RMSE) are given for the pre- and post-recalibration curves.

process in addition to NA12878 herself. Increasing the number of samples in the cohort rapidly improved both the sensitivity and specificity of the call set. As evidence mounts with more samples that a particular site is polymorphic, our confidence in the call increases and the site is more likely to be called (Fig. 6a).

Distinguishing true positive variants from sequencing and data processing artifacts is more difficult with few samples and, consequently, low aggregated coverage; adding more reads allows the error covariates to identify sites as errors using the variant recalibrator (Fig. 6b,c).

The combination of multi-sample SNP calling, variant quality recalibration using error covariates and imputation allows one to achieve

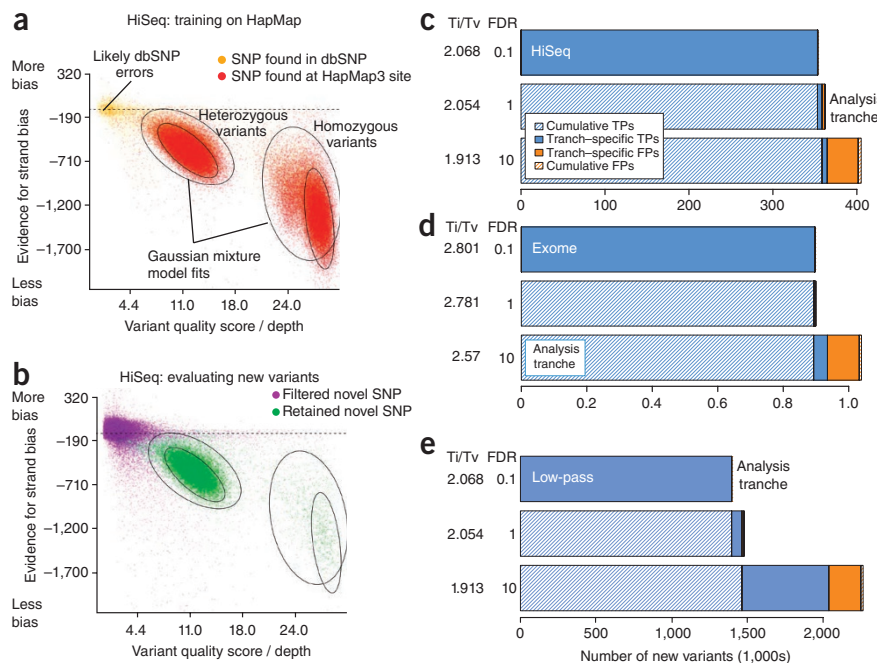
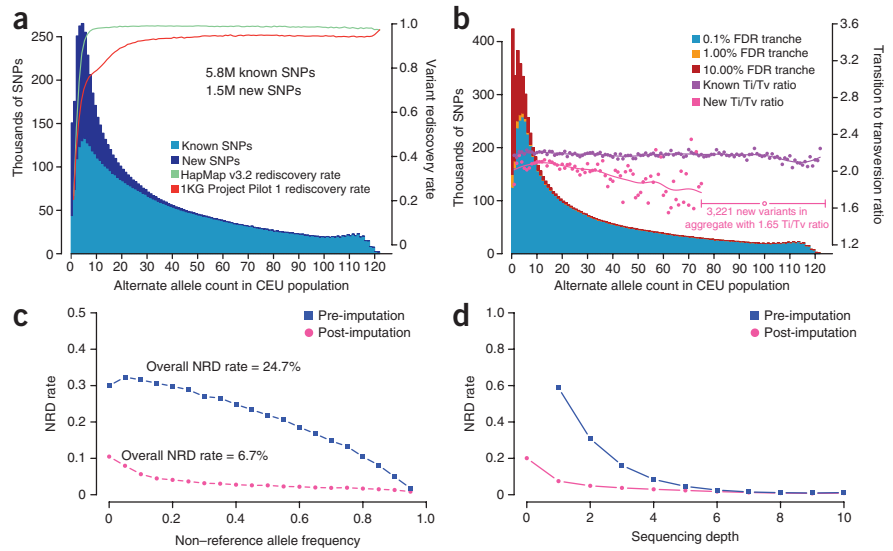


Figure 4 Results of variant quality recalibration on HiSeq, exome and low-pass data sets. (a) Relationship in the HiSeq call set between strand bias and quality by depth for genomic locations in HapMap3 (red) and dbSNP (orange) used for training the variant quality score recalibrator (left), (b) and the same annotations applied to differentiate likely true positive (green) from false positive (purple) new SNPs. (c–e) Quality tranches in the recalibrated HiSeq (c), exome (d) and low-pass CEU (e) calls beginning with (top) the highest quality but smallest call set with an estimated false positive rate among new SNP calls of <1/1000 to a more comprehensive call set (bottom) that includes effectively all true positives in the raw call set along with more false positive calls for a cumulative false positive rate of 10%. Each successive call set contains within it the previous tranche's true- and false-positive calls (shaded bars) as well as tranche-specific calls (solid bars). The tranche selected for further analyses here is indicated.

Figure 5 Variation discovered among 60 individuals from the CEPH population from the 1000 Genomes Project pilot phase plus low-pass NA12878. (a) Discovered SNPs by non-reference allele count in the 61 CEPH cohort, colored by known (light blue) and new (dark blue) variation, along with non-reference sensitivity to CEU HapMap3 and 1000 Genomes Project low-pass variants. (b) Quality and certainty of discovered SNPs by non-reference allele count. The histogram depicts the certainty of called variation broken out into 0.1%, 1% and 10% new FDR tranches. The Ti/Tv ratio is shown for known and new variation for each allele count, aggregating the new calls with allele count >74 because of their limited numbers. (c,d) Genotyping accuracy for NA12878 from reads alone (blue squares) and following genotype-likelihood based imputation (pink circles) called in the 61 sample call set as assessed by the NRD rate to HiSeq genotypes as a function of allele count (c) and sequencing depth (d).



a high-quality call set, both in aggregate and per sample, with very little data. The aggregated 61-sample set at 4× coverage includes only four times as much sequencing data as the HiSeq data, yet we discovered 3.2 million polymorphic sites in NA12878, which includes 97%, 91% and 87% of the variants in the HapMap3, 1000 Genomes Project Trio and HiSeq call sets, respectively, while also finding ~5 million additional variants among the 60 other samples.

Hard filtering versus variant quality score recalibration

Supplementary Table 3 lists the quality of call sets derived using our previous filtering approaches on all three datasets relative to the adaptive recalibrator described here. In all cases, the adaptive approach outperformed the manually optimized hard filtering previously developed for this calling system for the 1000 Genomes Project pilot data. This highlights two important points: first, that a principled integration of all covariates (which may have a complex correlation structure) should and does outperform single manually defined thresholds on covariates independently, with the added benefit of not requiring human intervention; and second, that an accurate ranking of discovered putative variants by

Supplementary Table 3 lists the quality of call sets derived using our previous filtering approaches on all three datasets relative to the adaptive recalibrator described here. In all cases, the adaptive approach outperformed the manually optimized hard filtering previously developed for this calling system for the 1000 Genomes Project pilot data. This highlights two important points: first, that a principled integration of all covariates (which may have a complex correlation structure) should and does outperform single manually defined thresholds on covariates independently, with the added benefit of not requiring human intervention; and second, that an accurate ranking of discovered putative variants by

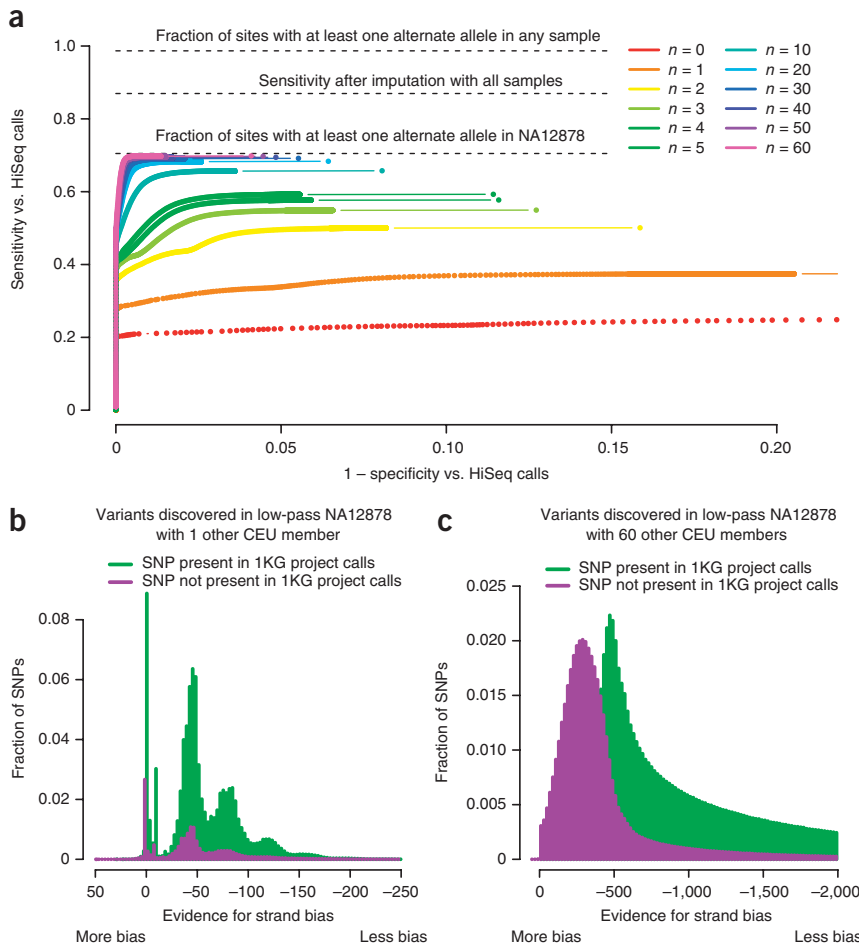


Figure 6 Sensitivity and specificity of multi-sample discovery of variation in NA12878 with increasing cohort size for low-pass NA12878 read sets processed with *N* additional CEPH samples. (a) Receiver operating characteristic (ROC) curves for SNP calls relating specificity and sensitivity to discover non-reference sites from the NA12878 HiSeq call set. The maximum callable sensitivity, 66%, is the percent of sites from the HiSeq call set where at least one read carries the alternate allele in the low-pass data for NA12878; it reflects both differences in the sequencing technologies (36–76-bp GAII for the low-pass NA12878 sample compared to 101-bp HiSeq) as well as the vagaries of sampling at 4× coverage. Because most of these missed sites are common and are consequently called in the other samples, imputation recovers ~50% of these sites. (b,c) Increasing power to identify strand-biased, likely false positive SNP calls with additional samples. Histograms of the strand bias annotation at raw variant calls discovered in the low-pass CEU data using NA12878 at 4× combined with one other CEU individual (b) and with 60 other individuals (c) stratified into sites present (green) and not (purple) in the 1000 Genomes Project CEU trio.



the probability that each represents a true site permits the definition of tranches for specificity or sensitivity (Fig. 4c–e) as appropriate to the needs of the specific project. Although the most permissive tranche includes almost all sites that have any chance of being true polymorphisms—critical for projects looking for single large-effect mutations—the vast majority of true polymorphisms are present in the highest quality tranche of data (data not shown).

Comparison of this calling pipeline to Crossbow

To calibrate the additional value of the tools described here, we contrasted our results with SNPs called on our raw NA12878 exome data using Crossbow²⁸, a package combining bowtie, a gapless read mapping tool based on the Burrows–Wheeler transformation²⁹, and SoapSNP for SNP detection¹⁵. We chose to perform this analysis on the exome data because its wide range of read depths and complex error modes make SNP calling a challenge, especially given the small number of new variants (~1,000 per sample) expected in this 28-Mb target. In **Supplementary Table 4**, the high-level results of the GATK and Crossbow calling pipelines are compared and contrasted. Key metrics such as the number of new SNP calls, their Ti/Tv ratio, the number of calls not seen in either the 1000 Genomes Project trio or the HiSeq data and the high nonsense and read-through rates indicate that the Crossbow call set has lower specificity than the GATK pipeline. This was true even after we applied an aggressive *P* value threshold (*P* < 0.01) for the base quality rank sum test¹⁵ to filter false-positive variants, which reduced the sensitivity of the HM3, 1000 Genomes Project and the HiSeq call sets by >3%. The intersection set between GATK and Crossbow is more specific but less sensitive than the calls unique to each pipeline (**Table 1**), a clear sign that despite the advances presented here, a lot of work remains to be done in perfecting calling in datasets like single sample exome capture. Although the value of the data processing and error modeling presented here is also clear, applying local realignment and base quality score recalibration (using publicly available, easy-to-use modules in the GATK) are likely to improve the results of the Crossbow pipeline.

DISCUSSION

The inaccuracy and covariation patterns differ strikingly between sequencing technologies (Fig. 3), which, if uncorrected, can propagate into downstream analyses. Accurately recalibrated base quality scores eliminate these sequencer-specific biases (Fig. 3) and enable integration of data generated from multiple systems. Although developed for early NGS datasets like those from the 1000 Genomes Project pilot, the impact of recalibration is still substantial even for data emerging today on newer sequencers like the HiSeq (2000). Together with local realignment, these two data processing methods eliminated millions of mostly false positive variants while preserving nearly all true variable sites, such as those in HapMap3 and 1KG Trio (**Table 2**). In single sample datasets, such as HiSeq and exome, without realignment and recalibration, these false variants account for more than a fifth of all of the new calls.

Even with very deep coverage, the naïve Bayesian model for SNP calling results in an initial call set with a surprisingly large number of false-positive calls. Although we expected 3.3 million known and 330,000 new non-reference sites in a single European sample sequenced genome wide, the initial HiSeq call set contains 3.5 million known and 800,000 new calls. The excessive number of variable sites, and the low Ti/Tv ratio in particular among the new calls, implies that ~600,000 of these variants are likely errors resulting from stochastic and systemic sequencing and alignment errors. The same calculations suggest that a similar fraction of the initial exome calls are likely false positives, and

more than 80% of the initial new low-pass SNP calls are likely errors. The adaptive error modeling developed here enabled us to identify these false-positive variants based on their dissimilarity to known variants, despite error rates of 50–80% among the new variants.

In each step of the pipeline, the improvements derive from the correction of systematic errors made in base calling or read mapping. By characterizing the specific NGS machine error processes and capturing our certainty, or lack thereof, that a putative variant is truly present in the sample or population, we delivered not a single concrete call set but a continuum from confident to less reliable variant calls for use as appropriate to the specific needs of downstream analysis. Mendelian disease projects can select a more sensitive set of calls with a higher error rate to avoid missing that single, high-impact variant, whereas community resource projects like the 1000 Genomes Project can place a high premium on specificity.

The division between SNP discovery and preliminary genotyping and genotype refinement (columns 2 and 3 of Fig. 1) avoids embedding in the discovery phase assumptions about population structure, sample relationships and the LD relationships between variants. Consequently, our calling approach applies equally well to population samples in Hardy–Weinberg equilibrium like mother–father–child trios or interbreeding families suffering from Mendelian disorders. Critically, our framework produces highly sensitive and specific variation calls without the use of LD and so can be applied in situations where LD information is unavailable or weak (many organisms) or would confound analytic goals such as studying LD patterns themselves or comparing Neanderthals and modern humans³⁰. Where appropriate, however, imputation can be applied to great value, as we demonstrated in the 61-sample CEU low-pass call set.

The analysis results presented here clearly indicate that even with our best current approaches we are still far from obtaining a complete and accurate picture of genetic variation of all types in even a single sample. Even with the HiSeq 10-bp paired-end reads, nearly 4% (~100 Mb) of the potentially callable genome is considered poorly mapped (**Supplementary Note**), and analysis of variants within these regions requires care. Nearly two thirds of the differences between the HiSeq and exome call sets can be attributed to different read mappings between BWA and MAQ.

The challenge of obtaining accurate variant calls from NGS data is substantial. We have developed an analysis framework for NGS data that achieves consistent and accurate results from a wide array of experimental design options including diverse sequencing machinery and distinct sequencing approaches. We have introduced here an integrated approach to data processing and variation discovery from NGS data that is designed to meet these specifications. Using data generated both at the Broad Institute and throughout the 1000 Genomes Project, we have shown that the introduction of improved calibration of base quality scores, local realignment to accommodate indels, the simultaneous evaluation of multiple samples from a population, and finally, an assessment of the likelihood that an identified variable site is a true biological DNA variant greatly improves the sensitivity and specificity of variant discovery from NGS data. The impending arrival of yet more NGS technologies makes even more important modular, extensible frameworks like ours that produce high-quality variant and genotype calls despite distinct error modes of multiple technologies for many experimental designs.

METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.

Note: Supplementary information is available on the Nature Genetics website.

ACKNOWLEDGMENTS

Many thanks to our colleagues in Medical and Population Genetics and Cancer Informatics and the 1000 Genomes Project who encouraged and supported us during the development of the Genome Analysis Toolkit and associated tools. This work was supported by grants from the National Human Genome Research Institute, including the Large Scale Sequencing and Analysis of Genomes grant (54 HG003067) and the Joint SNP and CNV calling in 1000 Genomes sequence data grant (U01 HG005208). We would also like to thank our excellent anonymous reviewers for their thoughtful comments.

AUTHOR CONTRIBUTIONS

M.A.D., E.B., R.P., K.V.G., J.R.M., C.H., A.A.P., G.d.A., M.A.R., T.J.F., A.Y.S. and K.C. conceived of, implemented and performed analytic approaches. M.A.D., E.B., R.P., K.V.G., G.d.A., A.M.K. and M.J.D. wrote the manuscript. M.A.D., M.H. and A.M. developed Picard and GATK infrastructure underlying the tools implemented here. M.A.D., S.B.G., D.A. and M.J.D. lead the team.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturegenetics/>.
Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. The 1000 Genomes Project Consortium. *et al.* A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
2. Yi, X. *et al.* Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* **329**, 75–78 (2010).
3. Ng, S.B. *et al.* Exome sequencing identifies the cause of a mendelian disorder. *Nat. Genet.* **42**, 30–35 (2009).
4. Lee, W. *et al.* The mutation spectrum revealed by paired genome sequences from a lung cancer patient. *Nature* **465**, 473–477 (2010).
5. Pleasance, E.D. *et al.* A comprehensive catalogue of somatic mutations from a human cancer genome. *Nature* **463**, 191–196 (2009).
6. Beroukhi, R. *et al.* The landscape of somatic copy-number alteration across human cancers. *Nature* **463**, 899–905 (2010).
7. Roach, J.C. *et al.* Analysis of genetic inheritance in a family quartet by whole-genome sequencing. *Science* **328**, 636–639 (2010).
8. Li, R. *et al.* SOAP2: an improved ultrafast tool for short read alignment. *Bioinformatics* **25**, 1966–1967 (2009).
9. Li, H., Ruan, J. & Durbin, R. Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Res.* **18**, 1851–1858 (2008).
10. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
11. Ning, Z., Cox, A.J. & Mullikin, J.C. SSAHA: a fast search method for large DNA databases. *Genome Res.* **11**, 1725–1729 (2001).
12. Ewing, B. & Green, P. Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Res.* **8**, 186–194 (1998).
13. Brockman, W. *et al.* Quality scores and SNP detection in sequencing-by-synthesis systems. *Genome Res.* **18**, 763–770 (2008).
14. Li, M., Nordborg, M. & Li, L.M. Adjust quality scores from alignment and improve sequencing accuracy. *Nucleic Acids Res.* **32**, 5183–5191 (2004).
15. Li, R. *et al.* SNP detection for massively parallel whole-genome resequencing. *Genome Res.* **19**, 1124–1132 (2009).
16. Drmanac, R. *et al.* Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays. *Science* **327**, 78–81 (2010).
17. Bentley, D.R. *et al.* Accurate whole human genome sequencing using reversible terminator chemistry. *Nature* **456**, 53–59 (2008).
18. Koboldt, D., Chen, K., Wylie, T. & Larson, D. VarScan: variant detection in massively parallel sequencing of individual and pooled samples. *Bioinformatics* **25**, 2283–2285 (2009).
19. Wheeler, D.A. *et al.* The complete genome of an individual by massively parallel DNA sequencing. *Nature* **452**, 872–876 (2008).
20. Mokry, M. *et al.* Accurate SNP and mutation detection by targeted custom microarray-based genomic enrichment of short-fragment sequencing libraries. *Nucleic Acids Res.* **38**, e116 (2010).
21. Shen, Y. *et al.* A SNP discovery method to assess variant allele probability from next-generation resequencing data. *Genome Res.* **20**, 273–280 (2010).
22. Hoberman, R. *et al.* A probabilistic approach for SNP discovery in high-throughput human resequencing data. *Genome Res.* **19**, 1542–1552 (2009).
23. Malhis, N. & Jones, S. High quality SNP calling using Illumina data at shallow coverage. *Bioinformatics* **26**, 1029 (2010).
24. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
25. Handsaker, R.E., Korn, J.M., Nemes, J. & McCarroll, S.A. Discovery and genotyping of genome structural polymorphism by sequencing on a population scale. *Nat. Genet.* **43**, 269–276 (2011).
26. McKenna, A.H. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
27. Browning, B.L. & Yu, Z. Simultaneous genotype calling and haplotype phasing improves genotype accuracy and reduces false-positive associations for genome-wide association studies. *Am. J. Hum. Genet.* **85**, 847–861 (2009).
28. Langmead, B., Schatz, M.C., Lin, J., Pop, M. & Salzberg, S.L. Searching for SNPs with cloud computing. *Genome Biol.* **10**, R134 (2009).
29. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
30. Green, R.E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010).
31. Gnirke, A. *et al.* Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nat. Biotechnol.* **27**, 182–189 (2009).
32. Ng, S., Turner, E., Robertson, P. & Flygare, S. Targeted capture and massively parallel sequencing of 12 human exomes. *Nature* **461**, 272–276 (2009).
33. Mckernan, K.J. *et al.* Sequence and structural variation in a human genome uncovered by short-read, massively parallel ligation sequencing using two-base encoding. *Genome Res.* **19**, 1527–1541 (2009).

© 2011 Nature America, Inc. All rights reserved.



ONLINE METHODS

Evaluating the quality of SNP calls. *Number of SNP calls and allele frequency.* The number of calls and frequency for multi-sample calling should follow relatively closely the neutral expectation for N individuals for small N :

$$\text{Number of polymorphic sites} \approx L \times \theta \sum_{i=1}^{2N} 1/i$$

where L is the number of confidently called bases and θ is the population-specific heterozygosity, genome wide of $\sim 0.8 \times 10^{-3}$ for CEPH individuals (H. Li, unpublished data). A surplus of variants, especially heterozygous variants for single samples or lower-frequency variants for populations, is a strong indicator of false positives.

dbSNP rate. Most variants are already catalogued in the dbSNP database of human variation. For a single European sample, $\sim 90\%$ of their true variants will appear in dbSNP build 129 (**Supplementary Table 5**), which will reach $\sim 99\%$ following the completion of the 1000 Genomes Project (**Supplementary Fig. 1**). For population-level SNP calls, the aggregate dbSNP rate for the call set decreases as more rare variants are found, which are less frequently found in dbSNP. Nevertheless, the per sample dbSNP rate should remain consistent across individuals. Note that presence in dbSNP is not an absolute confirmation that a variant is true (for example, see **Fig. 2** and **Fig. 4**), but because dbSNP build 129 contains 11.6 million SNP entries (only 0.4% of all genomic positions), relative differences between call sets with high dbSNP rates can be reasonably interpreted as quality differences.

Non-reference sensitivity and non-reference discrepancy (NRD) rate. For single samples, comparison with non-reference genotype calls from microarray chips, such as HapMap3 (~ 1.3 – 1.5 million sites), provides a good initial assessment of variant discovery sensitivity. With sufficient coverage, $>99\%$ of non-reference sites can generally be discovered. The NRD rate reports the percent of discordant genotype calls at commonly called non-reference sites on the chip and should reach $<1\%$ with sufficient coverage. Mathematical definitions of these terms are:

$$\text{NR sensitivity}(E, C) = \frac{|Enr \cap Cnr|}{|Cnr|}$$

$$\text{NRD rate}(E, C) = \frac{|\{i \in Enr \cup Cnr : Ei \neq Ci\}|}{|Enr \cup Cnr|}$$

X_i = Number of non-reference alleles for genotype call i in call set X

$$X_{nr} = \{i \in X : X_i > 0\}$$

E = Call set to be evaluated

C = Call set to be compared to

Transition/transversion ratio (Ti/Tv). The Ti/Tv ratio is a critical metric for assessing the specificity of new SNP calls. Inter-species comparisons³⁴ and previous sequencing projects (**Supplementary Table 6**) agree on a Ti/Tv ratio of ~ 2.0 – 2.1 for genome-wide datasets and 3.0 – 3.3 for exonic variation³⁵. The expected values for the Ti/Tv for known and new variants genome wide are 2.10 and 2.07, respectively, and in the exome target are 3.5 and 3.0, respectively. Currently the lower Ti/Tv ratio at new sites than at known sites is because of a combination of residual false positives lowering the Ti/Tv, a relative deficit of transitions due to sequencing context bias, as well as an apparently higher transition ratio at lower frequency variation. These uncertainties should limit the interpretation of minor differences in Ti/Tv ratios (<0.05), especially across sequencing technologies and datasets.

The Ti/Tv ratio for randomly assigned ‘variation’, such as results from systematic sequencing errors, alignment artifacts and data processing failures will be ~ 0.5 , as there are two transversion mutations for each transition. Given an expected Ti/Tv ratio, as above, and an observed Ti/Tv ratio from a call

set, an estimate of the fraction of false positive variants in the call set can be obtained by:

$$FDR_{\text{test}} = \frac{TiTv_{\text{observed}} - 0.5}{TiTv_{\text{expected}} - 0.5}$$

which should be bounded above by 100% (because of Ti/Tv ratios below 0.5) and a minimum false-positive rate (here assumed to be 0.1%) when the observed Ti/Tv exceeds the expected value.

Local multiple sequence realignment. We developed a local realignment algorithm that provides a consistent alignment among all reads spanning an indel. The algorithm begins by first identifying regions for realignment where (i) at least one read contains an indel, (ii) there exists a cluster of mismatching bases or (iii) an already known indel segregates at the site (for example, from dbSNP). At each region, haplotypes are constructed from the reference sequence by incorporating any known indels at the site, indels in reads spanning the site or from Smith-Waterman³⁶ alignment of all reads that do not perfectly match the reference sequence. For each haplotype H_i , reads are aligned without gaps to H_i and scored according to:

$$L(R_j | H_i) = \prod_k L(R_{j,k} | H_{j,i})$$

$$L(R_{j,k} | H_{j,i}) = \begin{cases} 1 - \epsilon_{j,k} & \approx 1 \quad R_{j,k} = H_{j,i} \\ \epsilon_{j,k} & R_{j,k} \neq H_{j,i} \end{cases}$$

$$L(H_i) = \prod_j L(R_j | H_i)$$

where R_j is the j th read, k is the offset in the gapless alignment of R_j and H_i and $\epsilon_{j,k}$ is the error rate corresponding to the declared quality score for the k th base of read R_j . The haplotype H_i that maximizes $L(H_i)$ is selected as the best alternative haplotype. Next, all reads are realigned against just the best haplotype H_i and the reference (H_0), and each read R_j is assigned to H_i or H_0 depending on whichever maximizes $L(R_j | H)$. The reads are realigned if the log odds ratio of the two-haplotype model is better than the single reference haplotype by at least five log units:

$$\frac{L(H_0, H_i)}{L(H_0)} = \frac{\prod_j \max[L(R_j | H_i), L(R_j | H_0)]}{\prod_j L(R_j | H_0)}$$

This discretization reflects a tradeoff between accuracy and efficient calculation of the full statistical quantities. Note that this algorithm operates on all reads across all individuals simultaneously, which ensures consistency in the inferred haplotypes among all individuals, a critical property for reliable indel calling and contrastive analyses such as somatic SNP and indel calling. The realigned reads are written to a SAM/BAM file for further analysis. The reads around a homozygous deletion, before and after local realignment, for Genome Analyzer reads from the 1000 Genomes Project and HiSeq, are shown in **Figure 2**.

Base quality score recalibration. We developed a base quality recalibration algorithm that provides empirically accurate base quality scores for each base in every read while also correcting for error covariates like machine cycle and dinucleotide context, as well as supporting platform-specific error covariates like color-space mismatches for SOLiD and flow-cycles for 454 (refs. 13–15,37,38). For each lane, the algorithm first tabulates empirical mismatches to the reference at all loci not known to vary in the population (dbSNP build 129), categorizing the bases by their reported quality score (R), their machine cycle in the read (C) and their dinucleotide context (D). For each category we estimate the empirical quality score:

$$\text{mismatch}(R, C, D) = \sum_{r \in R} \sum_{c \in C} \sum_{d \in D} \sum_{b, r, c, d \neq \text{bref}} br, c, d$$

$$\text{bases}(R, C, D) = \sum_{r \in R} \sum_{c \in C} \sum_{d \in D} |\{b, r, c, d\}|$$

$$Q_{\text{empirical}}(R, C, D) = (\text{mismatch}(R, C, D) + 1) / (\text{bases}(R, C, D) + 1)$$

These covariates are then broken into linearly separable error estimates and the recalibrated quality score Q_{recal} is calculated as:

$$\begin{aligned} \text{recal}(r, c, d) &= Q_r + \Delta Q(r) + \Delta\Delta Q(r, c) + \Delta\Delta Q(r, d) \\ \Delta Q &= Q_{\text{empirical}}(R, C, D) - \left(\sum_{\rho} \varepsilon_r \times N_r \right) / \text{bases}(R, C, D) \\ \Delta Q(r) &= Q_{\text{empirical}}(r, C, D) - Q_r - \Delta Q \\ \Delta Q(r, c) &= Q_{\text{empirical}}(r, c, D) - (\Delta Q_r + \Delta Q(r)) \\ \Delta Q(r, d) &= Q_{\text{empirical}}(r, C, d) - (\Delta Q_r + \Delta Q(r)) \end{aligned}$$

where each ΔQ and $\Delta\Delta Q$ are the residual differences between empirical mismatch rates and that implied by the reported quality score for all observations conditioning only on Q_r , or on both the covariate and Q_r ; Q_r is the base's reported quality score and ε_r is its expected error rate; $b_{r,c,d}$ is a base with specific covariate values, and r, c, d and R, C, D are the sets of all values of reported quality scores, machine cycles and dinucleotide contexts, respectively. The quality score and covariate distributions for four datasets before and after quality score recalibration are shown in **Figure 3**.

Multi-sample SNP calling. We apply a Bayesian algorithm for variant discovery and genotyping that simultaneously estimates the probability that two alleles A, the reference allele, and B, the alternative allele, are segregating in a sample of N individuals and the likelihoods for each of the AA, AB and BB genotypes for each of individual. Given D_i aligned bases at a specific genomic position for individual i , we estimate the genotype likelihoods GT_i of observing the D_i bases for each of AA, AB and BB genotypes according to the following equation:

$$\begin{aligned} \Pr\{D_i | GT_i\} &= \prod_j \Pr\{D_{i,j} | GT_i\} \\ \Pr\{D_i | GT_i = AB\} &= (\Pr\{D_{i,j} | A\} + \Pr\{D_{i,j} | B\}) / 2 \\ \Pr\{D_{i,j} | B\} &= \begin{cases} 1 - \varepsilon_{i,j} \\ \varepsilon_{i,j} \cdot \Pr\{B \text{ is true} | D_{i,j} \text{ is miscalled}\} \end{cases} D_{i,j} = B, \text{ otherwise.} \end{aligned}$$

where $\Pr\{D_{i,j} | GT_i\}$ is the probability of observing base $D_{i,j}$ under the hypothesized genotype GT_i ; $\Pr\{D_{i,j} | B\}$ and $\Pr\{D_{i,j} | A\}$ are the probability of observing base $D_{i,j}$ given that the true base is B or A, respectively; $\varepsilon_{i,j}$ is the probability of a base miscall given the quality score of base $D_{i,j}$; and $\Pr\{B \text{ is true} | D_{i,j} \text{ is miscalled}\}$ is the probability of B_{true} being the true chromosomal base given that b is a miscall (**Supplementary Table 7**). As these are raw likelihoods, no prior probabilities are applied.

Let us define $q_i = \{0,1,2\}$ as the number of alternate B alleles carried by individual i , so that $q = \sum_i q_i$ is the number of chromosomes carrying the B allele among all individuals. We estimate the probability that $q = X$ as:

$$\begin{aligned} \Pr\{q = X | D\} &= \frac{\Pr\{q = X\} \Pr\{D | q = X\}}{\sum_Y \Pr\{D | q = Y\}} \\ \Pr\{q = X\} &= \begin{cases} \theta / X \\ 1 - \theta \sum_{i=1}^{2N} 1/i \end{cases} X > 0 \text{ otherwise.} \\ \Pr\{D | q = X\} &= \sum_{GT \in \Gamma} \prod_i \Pr\{D_i | GT_i\} \\ \Gamma &= \left\{ GT \text{ where } \sum_i q_i = X \right\} \end{aligned}$$

where Γ is the set of all genotype assignments for the N individuals that contain exactly $q = X$ B alleles, $\Pr\{q = X\}$ is the infinite-sites neutral expectation

to observe X alternative alleles in $2N$ chromosomes with heterozygosity of θ , and GT_i and D_i are the i th individual's genotype and NGS reads, respectively. The sum over Γ involves potentially evaluating 3^N combinations but can be approximated by a heuristic algorithm like expectation-maximization through the introduction of a Hardy-Weinberg equilibrium assumption, using a greedy combinatorial search algorithm (**Supplementary Note**) or using an exact summation (H. Li, unpublished data). This algorithm emits the probability of a variant segregating at the site at some frequency:

$$\text{QUAL} = -10 \cdot \log_{10} [\Pr\{q = 0 | D\}]$$

represented conventionally by the Phred-scaled confidence, as well as the genotype assignments at the value that maximizes $\Pr\{q | D\}$. Only sites with $\text{QUAL} > Q50$ for deep coverage or $Q10$ for shallow coverage, respectively, are considered here as potentially variable sites.

Variant quality score recalibration. Given a set of putative variants along with SNP error covariate annotations, variant quality score recalibration employs a variational Bayes Gaussian mixture model (GMM)³⁹ to estimate the probability that each variant is a true polymorphism in the samples rather than a sequencer, alignment or data processing artifact. The set of variants $\{v_i\}$ are treated as an n -dimensional point cloud, each variant v_i positioned by its covariate annotation vector, \bar{v} . A mixture of Gaussians is fit to the set of likely true variants, here approximated by the variants already present in HapMap3 (**Fig. 4a**). Following training, this mixture model is used to estimate the probability of each variant call being true (**Fig. 4b**), capturing the intuition that variants with similar characteristics as previously known variants are likely to be real, whereas those with unusual characteristics are more likely to be machine or data processing artifacts.

Mathematically, we write the probability of a variant's vector of covariate values as the linear superposition of Gaussians:

$$\begin{aligned} \Pr\{v_i | GMM\} &= \sum_{k=1}^K \pi_k N(\bar{v}_i | \bar{\mu}_k, \Sigma_k) \\ \Pr\{\bar{\pi}\} &= \text{Dir}(\bar{\pi} | \alpha_0) = C(\alpha_0) \prod_{k=1}^K \pi_k^{\alpha_0 - 1} \\ \Pr\{\bar{\mu}, \Lambda\} &= N(\bar{\mu} | \bar{m}_0, (\beta_0 \Lambda_k)^{-1}) W(\Lambda_k | W_0, v_0) \end{aligned}$$

where K is the number of Gaussians in the mixture (GMM), and the last two equations are standard conjugate prior distributions over the parameters $\bar{\pi}$, $\bar{\mu}$ and Σ .

We then use an analog of the expectation-maximization algorithm³⁹ to learn the optimal parameters for the clusters using only variant calls at sites present in HapMap3. By restricting training to known polymorphic sites, the resulting GMM captures the distribution of covariate parameters for true SNPs. Consequently, we estimate the likelihood of each putative variant v_i being true under the learned GMM as:

$$\begin{aligned} L(v_i | GMM) &= \Pr\{v_i\} \Pr\{\bar{v}_i | GMM\} \\ L(v_i | GMM) &= (1 - FP_{\text{singleton}})^{AC} \Pr\{\text{novelty of } v_i\} \sum_{k=1}^K \pi_k N(\bar{v}_i | \bar{\mu}_k, \Sigma_k) \\ \Pr\{\text{novelty of } v_i\} &= \begin{cases} 97\% \text{ } v_i \text{ is in HapMap3,} \\ 37\% \text{ otherwise.} \end{cases} \end{aligned}$$

where $\Pr\{v_i\}$ is the prior expectation that the putative variant v_i is true, \bar{v}_i is the vector of covariate values for v_i , $FP_{\text{singleton}}$ is the false positive rate for singletons (50% here), and AC is the number of chromosomes estimated to carry the variant among all called samples. The prior probability of $\Pr\{v_i\}$ depends on whether it is present in HapMap3 and its frequency in the samples being called, given an estimate of the false positive rate for singletons. This model can be easily extended to include more training data, more prior information and/or more error covariates.

For convenience of presentation and analysis, we partition the raw SNP calls into tranches based on the Ti/Tv ratio of their new variants. For each desired new false discovery rate target (FDR_i), tranche_i is defined as:

$$\text{tranche}_i = \{SNP_j \text{ where } L(SNP_j | GMM) > T_i\}$$

$$T_i = \text{smallest } X \text{ where } \text{titv}(\{SNP_j \text{ is novel} \wedge L(SNP_j | GMM) > X\}) > TiTv_i$$

$$TiTv_i = FDR_i * (TiTv_{\text{expected}} - 0.5) + 0.5$$

The first tranche is exceedingly specific but less sensitive, and each subsequent tranche in turn introduces additional true positive calls along with a growing number of false positive calls. More specificity in the learned GMM translates into better-separated tranches, where all true variants have high likelihoods and appear in the lowest FDR tranches and all false ones have low likelihoods and are excluded. Downstream applications can select in a principled way more specific or more sensitive call sets or incorporate directly the

recalibrated quality scores to avoid entirely the need to analyze only a fixed subset of calls but rather weight individual variant calls by their probability of being real.

34. Ebersberger, I., Metzler, D., Schwarz, C. & Pääbo, S. Genomewide comparison of DNA sequences between humans and chimpanzees. *Am. J. Hum. Genet.* **70**, 1490–1497 (2002).
35. Freudenberg-Hua, Y. *et al.* Single nucleotide variation analysis in 65 candidate genes for CNS disorders in a representative sample of the European population. *Genome Res.* **13**, 2271–2276 (2003).
36. Durbin, R., Eddy, S., Krogh, A. & Mitchison, G. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids.* (Cambridge University Press, Cambridge, UK, 1998).
37. Dohm, J.C., Lottaz, C., Borodina, T. & Himmelbauer, H. Substantial biases in ultra-short read data sets from high-throughput DNA sequencing. *Nucleic Acids Res.* **36**, e105 (2008).
38. HUGO Consortium. *et al.* Mapping human genetic diversity in Asia. *Science* **326**, 1541–1545 (2009).
39. Bishop, C. *Pattern Recognition and Machine Learning* (Springer, New York, New York, USA, 2006).

Supplemental information for DePristo et al., “A framework for variation discovery and genotyping using next-generation DNA sequencing data”

Supplementary Figures and Tables

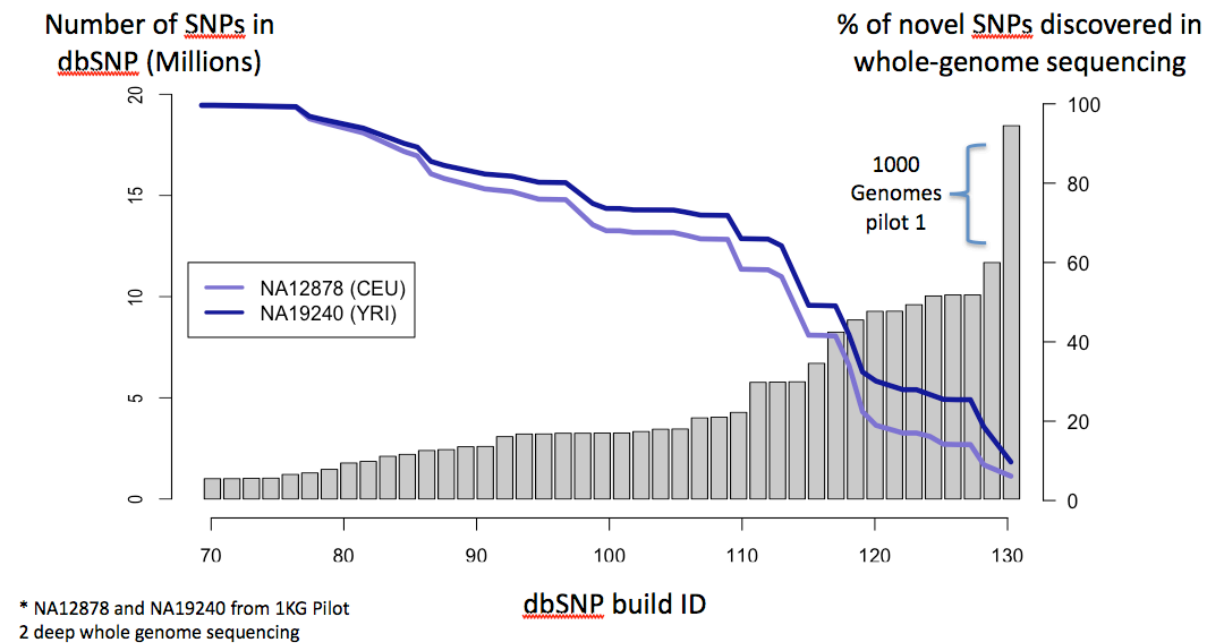


Figure S1: dbSNP rate for NA12878 and NA19240 from dbSNP build id 70 to 129 and then including additionally the CEU and YRI 1000 Genomes low-pass call sets

Table S1: Impact of local realignment in the HiSeq data set at sites in 1000 Genomes for which there is a homozygous indel for NA12878

Homozygous non-reference indel sites in NA12878 ^a	Homozygous indel sites with realigned reads	Spanning reads in raw alignments that do not contain the indel	Total # of realigned reads at homozygous non-reference indel sites ^b
124,568	116,259 (90% of all such sites)	1,448,627 (14.7% of all spanning reads)	1,083,125 (74.8% of reads without the indel)

a) According to 1000 Genomes CEU Trio indel call set

- b) 25% of the reads spanning indels do not contain the indel but match perfectly the reference sequence and so no realignment is necessary

Table S2: Number of regions, reads, and sites with significant mismatches affected by local realignment

Sequencing data set	Regions with realigned reads	Total size of affected regions (Mb)	Reads that were realigned	Sites with significant mismatch removal
HiSeq	947,765	21.3	6,621,462	1,749,840
Low-pass	2,627,318	57.3	16,586,650	991,532
Exome	49,170	1.1	149,449	106,182

Sites were considered to have a significant number of mismatches if over 15% of the bases (HiSeq and Exome data sets) or at least 2 bases (low-pass data) at the site mismatch the reference.

Table S3: Performance of hard-filtering and variant quality score recalibration

Call set	Site discovery						Comparison to NA12878 variants			
	No. of SNPs				Ti/Tv		HM3 concordance		1KG concordance	
	All	Known	Novel	dbSNP%	Known	Novel	NR sensitivity	NRD rate	NR sensitivity	NRD rate
HiSeq										
Recalibrated, MSA raw calls	4.18M	3.45M	722K	82.71	2.06	1.57	99.72	0.09	99.48	0.19
Hard filtered	3.53M	3.19M	351K	90.07	2.10	1.97	99.33	0.07	98.51	0.14
Variant recalibrated	3.58M	3.22M	362K	89.89	2.15	2.05	99.05	0.07	97.28	0.10
Low-pass*										
Recalibrated, MSA raw calls	13.4M	6.5M	6.9M	48.77	2.05	1.13	83.97	20.34	80.45	22.53
Direct Ti/Tv optimized	7.42M	5.05M	2.37M	67.98	2.10	1.82	82.79	20.23	78.14	22.17
Variant recalibrated call set	7.3M	5.8M	1.5M	79.7	2.18	2.05	83.02	20.26	76.99	22.01
Exome capture										
Recalibrated, MSA raw calls	18.5K	16.8K	1.7K	90.77	3.20	1.61	99.07	0.08	99.13	0.11
Hard filtered	15.9K	15.1K	807	94.93	3.38	2.74	97.24	0.08	97.10	0.11
Variant recalibrated	17.2K	16.2K	1039	93.96	3.27	2.57	98.49	0.08	98.38	0.11

* HM3 sensitivity and NRD rate only includes NA12878

Table S4: Comparison of NA12878 exome calls to reprocessing with Crossbow

	GATK	Crossbow	Intersection	Unique to GATK	Unique to Crossbow
No. known SNPs	16152	16086	15194	958	892
No. novel SNPs	1039	1806	799	240	1007
Known SNPs Ti/Tv	3.27	3.26	3.35	2.36	2.17
Novel SNPs Ti/Tv	2.57	1.94	2.77	2.04	1.50
HM3 NR sensitivity	98.5%	95.9%	95.3%	3.2%	0.6%
1000G Trio NR sensitivity	98.4%	95.4%	94.7%	3.7%	0.7%
Percent of calls not in 1000G trio	18.3%	23.9%	15.5%	55.8%	94.8%
HiSeq NR sensitivity	94.1%	90.9%	88.8%	5.3%	2.1%
Percent of calls not in HiSeq	3.5%	10.4%	2.1%	22.0%	80.5%
Percent synonymous variants	54.2%	52.8%	54.7%	48.0%	36.7%
Percent missense variants	45.5%	46.8%	45.1%	51.1%	61.6%
Percent nonsense/read-through	0.3%	0.4%	0.3%	0.9%	1.7%

Only includes calls in coding target regions.

Table S5: dbSNP 129 rates for several CEU and YRI samples using multiple sequencing technologies and SNP calling approaches

Sample	Population	dbSNP 129 rate	Sequencing technologies	SNP caller	Notes
NA12878	CEU	92%	Solexa, SOLiD and 454	Samtools + GATK	1000 Genomes official release
NA12891	CEU	92%	Solexa	Samtools + GATK	1000 Genomes official release
NA12892	CEU	92%	Solexa	Samtools + GATK	1000 Genomes official release
NA20431	CEU	90%	CG	CG	Complete genomics ¹
NA07022	CEU	90%	CG	CG	Complete genomics ¹

Table S6: Genome-wide and exome target Transition / Transversion (Ti/Tv) ratios expectations from published whole genome call sets against Human Genome build 36*.

Data set	Sequencing tech(s)	Year	SNP caller(s)	N sites	WGS		N sites	Exome	
					Known Ti/Tv	Novel Ti/Tv		Known Ti/Tv	Novel Ti/Tv
Initial resequencing projects									
Venter ¹	ABI	2007	Celera assembler	3.0M	2.10 ²	1.53 ²	15.2K	3.21 ²	2.54 ²
Watson ³	454	2008	Wheeler et al. caller	2.1M	2.13 ²	1.49 ²	12.2K	3.38 ²	1.90 ²
Single sample or trio NGS data sets									
Complete Genomics NA19240 ¹	CGI	2009	CGI	4.1M	2.14	2.09	20.2K	3.42	2.98
1000 Genomes CEU trio	Solexa, SOLiD and 454	2010	glfTrio, GATK	3.6M	2.08	2.02	17.6K	3.54	2.74
1000 Genomes YRI trio	Solexa, SOLiD and 454	2010	glfTrio, GATK	4.5M	2.09	2.07	25K	3.51	3.18
Weighted average					2.10	2.07	-	3.49	2.98
1000 Genomes low-coverage call sets									
CEU low-pass	Solexa, SOLiD and 454	2010	QCall, Mach, GATK	7.7M	2.10	1.90	45K	3.43	2.77
YRI low-pass	Solexa, SOLiD and 454	2010	QCall, Mach, GATK	10.6M	2.11	2.00	42K	3.55	2.91
Weighted average					2.10	1.96		3.48	2.84

- 1) Obtained from <http://huref.jcvi.org/>
- 2) Compared to dbSNP build 126, based on Wheeler et al. 2008
- 3) Obtained from <http://jimwatsonsequence.cshl.edu/>

Note that the transition / transversion ratio may depend on the properties of human genome reference build. These expected values should be calibrated for each major human genome reference version.

Table S7: Base miscalling confusion matrices by technology

Illumina (GA&HiSeq)				
	A	C	G	T
A	N/A	57.7%	17.1%	25.2%
C	34.9%	N/A	11.3%	53.9%
G	31.9%	5.1%	N/A	63.0%
T	45.8%	22.1%	32.0%	N/A
SOLiD				
	A	C	G	T
A	N/A	18.7%	42.5%	38.7%
C	27.0%	N/A	18.9%	54.1%
G	61.0%	15.7%	N/A	23.2%
T	40.5%	34.3%	25.2%	N/A
454				
	A	C	G	T
A	N/A	23.2%	42.6%	34.3%
C	19.7%	N/A	8.4%	71.9%
G	71.5%	6.6%	N/A	21.9%
T	43.8%	37.8%	18.5%	N/A

Supplementary Notes

Data generation

NA12878 HiSeq data

We sheared 1-3 ug of genomic DNA to a range of 100-700bp using the Covaris E210 instrument. DNA fragments were end-repaired and phosphorylated, followed by adenylation of 3' ends. Standard paired end adaptors were ligated according to the manufacturer's protocol (Illumina). We performed Qiagen min-elute column based cleanups between all enzymatic steps. Adapter ligated fragments were purified with preparatory gel electrophoresis (4% agarose, 85volts, 3 hours) and two bands were excised (500-520bp and 520-540bp) resulting in two libraries per sample with inserts averaging 380bp and 400bp respectively. DNA was extracted from gel bands using Qiagen min-elute columns. The entire volume of final purified fragments was enriched via PCR with Phusion polymerase for 10 cycles.

Libraries were quantified using a Sybr qPCR protocol with specific probes for the ends of the adapters. The qPCR assay measures the quantity of fragments properly adapter ligated that are appropriate for sequencing. Based on the qPCR quantification, libraries were normalized to 2nM and then denatured using 0.1 N NaOH. Cluster amplification of denatured templates occurred according to manufacturer's protocol (Illumina) using cBot reagent plates and HiSeq Flowcells (Illumina cat# PE-401-1001). Sybr Green dye was added to all flowcell lanes to provide a quality control checkpoint after cluster amplification to ensure optimal cluster densities on the flowcells. Flowcells were paired end sequenced with 101 bp reads on HiSeq2000s, using HiSeq Sequencing-by-Synthesis kits (Illumina cat# PE-401-1001) and analyzed with the Illumina v1.8 pipeline. Standard quality control metrics including error rates, % passing filter reads, and total Gb produced were used to characterize process performance prior to downstream analysis. The final data set includes 16 lanes totaling ~64x coverage of the genome.

NA12878 whole exome hybrid capture data

We sheared 1-3 ug of genomic DNA to a range of 100-300bp using the Covaris E210 instrument. DNA fragments were end-repaired and phosphorylated, followed by adenylation of 3' ends. Standard paired end adaptors were ligated according to the manufacturer's protocol (Illumina). We performed Ampure Bead-based cleanups between all enzymatic steps using the Bravo liquid handling platform (Agilent). The entire volume of final library fragments was enriched via PCR with Pfu (Agilent) polymerase for 6 cycles.

Libraries were quantified using an automated picogreen fluorescent assay compared against a standard curve of known samples and normalized to 25ng/ul prior to hybridization.

Hybridization of 500ng of library with 500ng of biotin-linked RNA (Agilent), designed specifically to the desired exome targets, was incubated at 65°C for 72 hours. Capture of resulting DNA-RNA duplexes was performed by the addition of streptavidin M280 beads (Invitrogen). Multiple washes at high stringency removed off-target material and any non-hybridized fragments. Desired fragments were PCR amplified directly off of beads using primers specific to the universal library sequences on the ends of captured fragments. The resulting captured libraries were quantified using picogreen in addition to a Sybr qPCR protocol (KAPA biosystems) with specific probes for the ends of the adapters.

Based on the qPCR quantification, libraries were normalized to 2nM and then denatured using 0.1 N NaOH. Cluster amplification of denatured templates occurred according to manufacturer's protocol (Illumina) using V2 Chemistry and V2 Flowcells (1.4mm channel width). Sybr Green dye was added to all flowcell lanes to provide a quality control checkpoint after cluster amplification to ensure optimal cluster densities on the flowcells. Flowcells were sequenced on Genome Analyzer II's, using V3 Sequencing-by-Synthesis kits and analyzed with the Illumina v1.3.4 pipeline. Standard quality control metrics including error rates, % passing filter reads, and total Gb produced were used to characterize process performance prior to downstream analysis.

NA12878 and 60 sample CEPH low-pass from 1000 Genomes

For the low-pass analysis, we used the publically available sequencing data from 61 individuals in the pilot phase of the 1000 Genomes Project. All 60 samples from the Caucasian (CEU) population of the low-pass wing were used; these individuals were sequenced to approximately 4x average coverage genome-wide on a variety of sequencing platforms: Illumina/GA, 454, and SOLiD. Additionally, we used the Illumina sequencing data of the CEU daughter (NA12878) from the high-pass trio wing of the 1000 Genomes pilot; as she was sequenced to over 30x coverage genome-wide, we downsampled her data to an average coverage of 4x. The downsampling was achieved by using only twenty-five read groups (lanes) from the high-pass data (ERR001751-ERR001775). The CEU low-pass data set includes 963 lanes of Solexa single- and paired-end reads, 95 lanes of SOLiD reads, and 906 lanes of 454 lanes comprising a total of ~600 Gb of sequence with on average 145x, 81x, and 42x coverage per platform, respectively.

Duplicate removal was performed using samtools rmdupse. An initial version of the GATK quality score recalibration tool was applied by the DCC to the GA and 454 BAMs; subsequent improvements to this tool include separately calibrating the first and second reads of a pair, reference-bias correction for SOLiD reads, as well as storing the original base qualities, which were not retained in the 1000 Genomes released BAMs and so could not be recalibrated here. Consequently, local multiple-sequence realignment was performed after quality score recalibration on all 61 individuals simultaneously.

Base miscalling confusion matrices

In order to account for biased miscalling for all three platforms, we tabulated reference and miscalled bases on chromosome 1 in the original 1000 Genomes CEU sample NA12878 sequenced with Illumina/GA, SOLiD, and 454 reads as part of the trio wing of the pilot project. Only loci containing exactly one non-reference base, at least 20x depth, where all bases had base qualities > Q20 and all reads had mapping qualities > Q30, and at least a Q50 homozygous reference genotype according to the genotyping algorithm caller presented below using the unified miscalling model (where $\Pr\{B \text{ true} \mid b \text{ miscall}\} = 1/3$) were considered. Such sites exhibit systematic miscalling biases of each instrument, from which we can calculate:

$$\Pr\{B_{true} \mid b_{miscalled}\} = \frac{\text{count}(b_{miscall} \mid B_{true})}{\sum_{X \neq B_{true}} \text{count}(b_{miscall} \mid X)}$$

These base miscalling confusion matrices for Solexa, SOLiD, and 454 are given in Table S3. Because they depend on particularities of the base calling algorithm applied for each technology which do change over time, these confusion matrices should be recalculated periodically.

Evaluating the quality of detected variation

We obtained HapMap3.2 consensus genotypes from hapmap.org, dbSNP 129 mapped to HG18 from the UCSC genome browser, and 1000 Genomes trio-aware SNP and indels calls from the 1000 Genomes DCC. Note that the 1000 Genomes CEU Trio call set is the intersection of SNP sites called by the GATK SNP caller (presented here, with hard-filtering as defined below) and a trio-aware extension to samtools

(<http://genome.sph.umich.edu/wiki/GlfTrio>) using three sequencing technologies to ~120x total depth. Genotypes for NA12878 and her parents were derived from the trio-aware caller and not the GATK caller.

Detecting reads from duplicate molecules

One variable and potentially large source of variation miscalls in NGS is non-independent sampling of DNA molecules during sequencing, such as occurs with repeated sampling of molecules that are molecular duplicates of one another. The PCR amplification steps involved in the majority of NGS library construction techniques can introduce significant biases due to preferential amplification of shorter molecules, molecules without extreme GC composition, etc. which will cause the sequencing to be a non-random sampling of the source genome. This is particularly problematic if any single molecule experiences a PCR error early in amplification as this error is propagated and sampled many times during sequencing.

To correct this problem we have developed an algorithm to detect and mark molecules that are probable duplicates of one another. This algorithm is simplified by the assumption that it is unlikely to sample the same exact molecule more than once from the source genome given true random sampling. Given current NGS protocols it is clear that this does in fact occur, but at an acceptably low rate:

De-duplication rate penalties by sequencing design strategy

Sequencing Application	Average #Molecules in Library	Read Length	Average #Molecules Sampled	Molecules Sampled > 1 times
30X Whole Genome	5bn	2x101b	~450m	4.4%*
4X Whole Genome	5bn	2x101b	~60m	0.6%
100X Whole Exome	500m	2x76b	~20m	2.0%

* Note that typically multiple independent libraries are created to see to such depth, thereby reducing the penalty of overmarking 'duplicate' molecules

Our duplicate marking algorithm relies on sequencing reads having been mapped to the genome to identify reads and read pairs that share the same start positions on the genome and mark these as duplicates.

Concretely, the steps taken to identify duplicate molecules are:

1. For each molecule, or cluster, sequenced identify the putative genomic position and strand for the 5'-most (with respect to the read) bases of each read originating from the molecule. The mapping positions reported by the aligner are then adjusted for any soft or hard clipping to determine the most

likely location of the 5'-most base whether or not that base has been mapped to the genome. The 5'-most bases are used as it is expected that these are the bases flanked by the sequencing adapters, which in turn are used as universal amplification sites during library construction.

2. Identify molecules where paired-reads have been performed and both reads have been mapped to the genome, and group these molecules by the genomic position and strand computed in step #1.
3. In groups of size > 1 , mark reads from all except one molecule as duplicates.
4. Identify molecules with only a single mapped read such that the mapped read's position and strand are identical to one end of a molecule with mapped paired-end reads and mark these molecules as duplicates also.
5. Group the remaining molecules with only a single mapped read by genomic position and strand and, similar to step #3, mark as duplicates all but one molecule in each group.

Within a group of duplicate molecules a simple heuristic is used to determine which molecule and hence which reads to retain. The base quality scores of each read are summed, ignoring those bases that are below Q15, and the read with the highest sum of quality scores is retained.

This algorithm is implemented in the program called MarkDuplicates in the Picard suite of tools (<http://picard.sourceforge.net/>). The program reads a SAM or BAM file as input and produces as output a SAM or BAM file with duplicate records retained but flagged.

Multi-sample SNP discovery and genotyping

The mathematical formulation of the multi-sample genotyping algorithm is given in Box 1. The genotyping algorithm makes a significant assumption by not directly modeling sequencing and alignment errors in the Bayesian framework, namely that any read present at a site actually belongs there. In actuality though, a high enough percentage of reads are misaligned as to affect the accuracy of the calling. To that end, we instituted several filters that must be passed in order for a given base to be used in the genotype likelihoods calculation: its phred-scaled base quality must be at least Q20, the mapping quality of its read must be at least 20, the read and its mate pair (for paired-end reads) must lie on the same chromosome, fewer than 10% of the bases on the read are permitted to mismatch the reference in a 20bp window on either side of the given base, and the read cannot be flagged as failing vendor quality filters or as originating from a duplicate molecule. Any base that passes all of these filters has an extremely high probability of actually originating at the given position.

The discovery algorithm is comprised of two phases: calculating the per-sample genotype likelihoods and then using a heuristic grid search to determine the most likely alternate allele frequency and genotype conformation over all samples. In the first phase, we use the bases at the position in question belonging to a given sample to calculate the likelihoods of the potential diploid genotypes for that sample (Box 1). As a heuristic improvement in the complexity of the problem we make the assumption that any given site is at most bi-allelic, choosing the most likely alternate allele based on the total sum of base qualities for each of the three possible non-reference bases. While tri-allelic sites, though rare, certainly do exist in diploid organisms, this assumption enables us to calculate only three likelihoods per sample (representing the possibility of the sample's being homozygous reference, heterozygous, or homozygous variant) instead of all ten possible diploid genotypes, which vastly improves the running time of the algorithm. We note that while this assumption may affect the genotype assignments at truly tri-allelic sites, it should not affect our ability to discover those sites.

Whereas the initial phase of the algorithm is run per sample, the second stage combines the genotype likelihoods over all samples in order to determine the most likely alternate allele frequency in the cohort. The likelihood for a given set of genotype assignments at a given frequency is simply the product of the genotype likelihoods for each sample given that sample's assigned genotype (Box 1). We then apply a population genetic prior to the allele frequency likelihoods based on θ , the population specific heterozygosity, and choose the most likely allele frequency and associated genotype assignments. The variant quality score for a polymorphic call is given as $-10 \cdot \log_{10}(\text{probability that the site is actually monomorphic})$.

In theory, for each possible alternate allele frequency, we would need to calculate the likelihood of each possible conformation of genotypes in the cohort under that frequency; however, as the number of possible genotype assignments is exponential in the number of samples, this calculation becomes intractable for larger cohorts. In practice, we have found that an excellent estimate of the most likely genotype assignment for a given allele frequency, f , is the most likely assignment for $f-1$ with another alternate allele added to one of the samples. By using a best-first search algorithm, calculating the likelihood of f is linear in the number of samples, as we can simply iterate over each sample and determine, using the per-sample genotype likelihoods, the most likely recipient of the new allele given the conformation at $f-1$. The fact that this greedy best-first algorithm converges once it hits a local maximum allows us to employ a further significant optimization: we can terminate the search for the most likely allele frequency early whenever the

likelihood for a given frequency is significantly lower than the maximum likelihood previously observed for any non-zero frequency.

We use indel calls in a post-process filtering step after SNP genotyping to remove those SNPs that are overly close to indels^{2,3}; these SNPs are usually an artifact of the partial discovery of an indel and are almost entirely false positives. We identify a site as potentially containing an indel if the indel occurs in a large enough fraction of the reads at a site (30% for high-pass data and 3% for low-pass) and is consistent within the reads (i.e. there is just a single consensus indel seen in the reads). Any SNPs called within 10bp on either end of these indels are then marked as filtered out of the callset. Note that these heuristic indel calls suffice to avoid SNP artifacts around misaligned indels, regardless of whether these indels result from machine artifacts or are truly segregating in the sample(s). Sites found to overlap the indel mask, and SNPs found in clusters (three or more within a ten-base-pair window) were filtered out of the call set³.

Variant Quality Score Recalibration

For the Dirichlet prior distribution over the mixing coefficients $\bar{\pi}$:

$$\Pr\{\bar{\pi}\} = \text{Dir}(\bar{\pi} \mid \alpha_0) = C(\alpha_0) \prod_{k=1}^K \pi_k^{\alpha_0 - 1}$$

where Dirichlet parameter α_0 is chosen to be 1×10^{-4} . For the Gaussian-Wishart prior over the mean and precision ($\Lambda^{-1} = \Sigma$) of each Gaussian in the mixture:

$$\Pr\{\bar{\mu}, \Lambda\} = \Pr\{\bar{\mu} \mid \Lambda\} \Pr\{\Lambda\} = N(\bar{\mu}_k \mid \bar{m}_0, (\beta_0 \Lambda_k)^{-1}) W(\Lambda_k \mid W_0, \nu_0)$$

we use a shrinkage parameter β_0 of 1×10^{-4} and ν_0 is given by the degrees of freedom which for this model is the number of covariates + 2 = 6.

Only a subset of known variants are used for clustering in order to avoid training on poorly determined annotations for variants with little sequencing data. In particular, only variants satisfying the following criteria are used to train the Gaussian mixture model:

Parameters for variant quality score recalibration

	HiSeq	Exome	Low-pass
Max. number of Gaussians to learn	16	8	6
Min. variant quality score for training	300	2800	1000
Max standard deviation from mean annotation value for inclusion	3.5	3.5	4.5
Max. percent of reads at a variant to be including in training	10%	10%	10%

Standard hard filters

For many projects, including our contributions to all three wings of the 1000 Genomes Project, we used a variety of hard filtering and optimization approaches to select high quality call sets from the raw calls. Here we list, for completeness, the application of our standard hard-filters for deep coverage and the predecessor to the variant recalibrator, a Ti/Tv-based standard-bias optimizer⁴, applied to the CEPH 61 sample set. Although it is possible to produce a reasonable quality call set using these approaches, the adaptive error modeling used by the variational Bayes recalibrator is able to better identify true positive variation with limited subjective intervention. In Table S9 we contrast the performance of the hard filtering approach to the variant recalibrator, and the improved specificity achievable with the recalibrator is clear. As with the variant recalibrator, even these hard filters become increasingly selective with additional samples, so that multiple deep read sets analyzed together provide a better trade-off in sensitivity and specificity.

For deep whole genomes, we filter out any SNPs matching the following criteria:

- The SNP cluster and proximity to indels filter as in the main analysis, or
- Greater than 10% of aligned reads at a site have mapping quality 0 (MAPQ0) among at least 40 reads, or
- $SB > -0.1$, or
- Depth of coverage above 120

For deep whole exomes, we filter out any SNPs matching the following criteria:

- The SNP cluster and proximity to indels filter as in the main analysis, or
- Greater than 10% of aligned reads at a site have mapping quality 0 (MAPQ0) among at least 40 reads, or
- $SB > -0.1$, or
- Quality over depth (QD) < 5 , or
- $H_{Run} > 3$

Additional annotation details

For SLOD: for each site, the procedure resulted in an estimate of allele frequencies in the population, an estimate of genotype likelihoods and posteriors for each member of the population, a LOD in favor of a site being variant, and an SB value measuring the strand bias in the non-reference allele. Under the null hypothesis, for a site detected as a variant the true non-reference allele frequency in the forward direction equals the true non-reference allele frequency in the reverse direction. Under the alternative hypothesis, where the site is not a variant but rather error prone, the non-reference allele frequency in the forward direction equals the estimated allele frequency and the non-reference allele frequency in the reverse direction equals to zero or vice-versa. SB is simply the log of the ratio of likelihood densities computed for the best supporting alternative hypothesis versus that of the null hypothesis.

The HaplotypeScore annotation associated with a SNP call at position POS is calculated by first (1) determining the two more prevalent 21 bp haplotypes around POS and then (2) calculating the probability of each read covering POS being sampled from either of these haplotypes:

1. Each read is enqueued into a priority queue with the priority being the sum of base quality scores within the 21 bp window (POS +/- 10 bp on each side). The set of perfect matching reads is set to ².
2. The read with the greatest sum of quality scores is taken from the queue. If it matches exactly any of the putative haplotypes, it is added to that haplotype's read set. Bases that are present in the new read but not in its haplotype are appended to the haplotype. If no exact match is available, a new haplotype is created with just this read in its read set. This process continues until the priority queue is exhausted and all reads have been placed into exact haplotype match sets.
3. We then construct consensus haplotypes with bases and quality scores for the two haplotypes with the greatest sum of base quality scores across all reads in their read sets.

4. Finally, the HaplotypeScore equation is determined for reads spanning POS against these two most common haplotypes.

The following figure depicts the intuition behind the Haplotype Score annotation:

Two segregating haplotypes

Read1 AAGC**T**CG
Read2 AAGC**A**CGA
Read3 AAGC**T**CGAT
Read4 AAGC**A**CGAT
Read5 AGC**A**CGAT
Read6 GC**T**CGAT

A/T polymorphic explains the reads well, so has a low haplotype score

Three segregating haplotypes

Read1 AA**G**C**T**CG
Read2 AA**G**C**A**CGA
Read3 AA**C**C**T**CGAT
Read4 AA**C**C**A**CGAT
Read5 A**G**C**A**CGAT
Read6 **C**C**T**CGAT

Inconsistent variation between A/T and C/G bases likely due to mapping artifacts, so has a high haplotype score

Likelihood-based genotype refinement with imputation

Beagle⁵ was used to refine genotypes in the low-pass 61-sample data set using likelihoods obtained during multi-sample SNP calling with default parameters and no external reference panel. In sites where no likelihoods were available for a sample due to lack of coverage, a model of uniform likelihoods was adopted. Genotypes were updated to those with the greatest posterior probabilities according to Beagle.

Specifically, given a set of variant sites obtained by the methods described in the preceding sections, the final step in variant discovery is the refinement and improvement of the obtained sample genotypes. In order to carry out this sample genotype improvement, the imputation software package Beagle 3.2 was used⁵ on the low-pass CEU population of 60 individuals, augmented by NA12878 downsampled to 4x coverage, after such data set was filtered using the Variant Recalibration procedure described above. For this application no reference panel was used, so that Beagle was only used to infer missing genotypes and to provide posterior genotype probabilities. The GATK was used to write the input to Beagle in its required format, using the genotype likelihoods from the SNP caller in Phase 2. In samples for which genotypes were missing (for example, at sites with no coverage), a uniform likelihood model was used (i.e. the likelihoods for a site being homozygous reference, heterozygous, or homozygous variant were set to an equal value of 1/3). After Beagle was run, the GATK was used again to parse the resulting output genotypes and posterior probabilities in order to produce an updated call set.

Two partitions of the data sets have proven useful for analyzing genotype imputation accuracy: a partition either by the non-reference allele frequency (AF) or the sequencing depth at each site. In both cases, NRD rates are computed for each bin. Figures 5(c) and 5(d) show the resulting NRD rates before and after imputation for each partition, as well as the global NRD rate from the whole data set as computed above. As Figure 5(d) illustrates, genotyping accuracy increases significantly at all points, and the increase is especially dramatic in low-depth sites. Even at sites where there was no sequencing coverage, we are able to recover genotypes with 20.9% NRD rate just by using imputation. With just one read, the NRD rate decreases from 58.8% to 7.6%. At the given target depth of 4x, the NRD rate improves from 8.5% to 3.0% with imputation. Importantly, the improvement in genotyping accuracy with 6 or more reads is marginal, which is the reason why imputation was only applied to the low-pass data and not to deep coverage data sets.

Comparison with Crossbow

To compare our results to existing data processing tools for next-generation DNA sequencing, we applied the Crossbow package ⁶ to the NA12878 whole exome sequence data. Specifically, the original machine output fastq files were aligned with bowtie ⁷ (-M 2, paired-end mode), and SNPs called with SoapSNP ⁸ (using the binomial probability calculation for greater accuracy). We considered only those sites identified as having QUAL score ≥ 20 and $P \geq 0.01$ for the rank sum test to determine if two alleles of a possible HET call have the same sequencing quality. The results of this comparison are summarized in Table S11.

Data availability

These data sets are available on the Genome Sequencing and Analysis website:

<http://www.broadinstitute.org/gsa/wiki/>

References

1. Drmanac, R. *et al.* Human Genome Sequencing Using Unchained Base Reads on Self-Assembling DNA Nanoarrays. in *Science* Vol. 327 78-81 (2010).
2. Bentley, D.R. *et al.* Accurate whole human genome sequencing using reversible terminator chemistry. in *Nature* Vol. 456 53-9 (2008).
3. Li, H., Ruan, J. & Durbin, R. Mapping short DNA sequencing reads and calling variants using mapping quality scores. in *Genome Research* Vol. 18 1851-1858 (2008).
4. The 1000 Genomes Project Consortium. A map of human genome variation from population scale sequencing. in *Nature* (2010).
5. Browning, B.L. & Yu, Z. Simultaneous genotype calling and haplotype phasing improves genotype accuracy and reduces false-positive associations for genome-wide association studies. in *Am J Hum Genet* Vol. 85 847-61 (2009).
6. Langmead, B., Schatz, M.C., Lin, J., Pop, M. & Salzberg, S.L. Searching for SNPs with cloud computing. in *Genome Biol* Vol. 10 R134 (2009).
7. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. in *Genome Biol* Vol. 10 R25 (2009).
8. Li, R. *et al.* SNP detection for massively parallel whole-genome resequencing. in *Genome Research* Vol. 19 1124-1132 (2009).