

Variation graph toolkit improves read mapping by representing genetic variation in the reference

Erik Garrison¹, Jouni Sirén¹, Adam M Novak² , Glenn Hickey², Jordan M Eizenga², Eric T Dawson^{1,3,4}, William Jones¹, Shilpa Garg⁵, Charles Markello², Michael F Lin⁶, Benedict Paten² & Richard Durbin^{1,4} 

Reference genomes guide our interpretation of DNA sequence data. However, conventional linear references represent only one version of each locus, ignoring variation in the population. Poor representation of an individual's genome sequence impacts read mapping and introduces bias. Variation graphs are bidirected DNA sequence graphs that compactly represent genetic variation across a population, including large-scale structural variation such as inversions and duplications¹. Previous graph genome software implementations^{2–4} have been limited by scalability or topological constraints. Here we present *vg*, a toolkit of computational methods for creating, manipulating, and using these structures as references at the scale of the human genome. *vg* provides an efficient approach to mapping reads onto arbitrary variation graphs using generalized compressed suffix arrays⁵, with improved accuracy over alignment to a linear reference, and effectively removing reference bias. These capabilities make using variation graphs as references for DNA sequencing practical at a gigabase scale, or at the topological complexity of *de novo* assemblies.

For small genomes, it is possible to study genetic variation by assembling whole genomes and then comparing them via whole-genome comparison^{6,7}. For large genomes, such as the human genome, complete and accurate *de novo* genome assembly is impractical because of repeat complexity and scale. Therefore, prior information is used to interpret new sequence data in their correct genomic context. The current practice is to align sequence reads to a single high-quality reference genome sequence that represents one haplotype at each location in the genome. Although this approach is much faster than *de novo* assembly, and simplifies discovery and reporting of genetic variants, it leads to mapping biases toward variants matching the reference sequence and away from alternative variants. There will even be some sequence in each new sample that is entirely absent in the reference⁸.

To avoid these biases, data would need to be aligned to a “personalized” reference sequence that already incorporates the individual's variants⁹, but in general it is not known what variants are present in a sample before aligning data from it. However, most differences between any one genome and the reference are segregating in the population. Thus, a reference structure that represents known shared

variation will contain most of the correct personalized sequence for any individual.

The natural computational structure for doing this is the sequence graph¹. Sequence graphs or equivalent structures have been used previously to represent multiple sequences that contain shared differences or ambiguities in a single structure. For example, multiple sequence alignments have a natural representation as partially ordered sequence graphs¹⁰. The variant call format¹¹ (VCF), which is a common data format for describing sets of genome sequences, can be understood as defining a partially ordered graph similar to those implied by a multiple sequence alignment. Related structures frequently used in genome assembly include the De Bruijn graph¹² and string graph¹³, which collapse long repeated sequences, so the same nodes are used for different regions of the genome. Graphs to represent genetic variation have previously been used for microbial genomes and localized regions of the human genome such as the major histocompatibility complex².

We define a variation graph as a sequence graph together with a set of paths representing possible sequences from a population (Fig. 1). Recently, software packages have been introduced that support a subset of variation graphs that reflect local variation away from a linear reference^{2,3}, formalizing approaches introduced in FreeBayes and the GATK HaplotypeCaller for the 1000 Genomes Project analysis^{14–16}. Our model goes beyond these in that it does not require the graph to be based on an initial linear reference, or indeed directionally ordered, and thus supports cycles and inversions. *vg* is the first openly available tool with these properties to scale practically to the multi-gigabase scale required for whole vertebrate genomes.

The core data model, data structures and algorithms, and implementation of *vg* are described in the Online Methods, with further details in the **Supplementary Note**. Indicative memory and compute runtime requirements are given in **Supplementary Table 1**. Below we present results demonstrating the mapping functionality of *vg*. Variant calling using *vg* against a variety of different human genome variation graphs is described elsewhere¹⁷.

For a species such as human, with only 0.1% nucleotide divergence on average between individual genome sequences, over 90% of 100-bp reads will derive from sequence exactly matching the reference. Therefore, new mappers should perform at least as well for linear

¹Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK. ²UC Santa Cruz Genomics Institute, University of California, Santa Cruz, California, USA. ³National Cancer Institute, Rockville, Maryland, USA. ⁴Department of Genetics, University of Cambridge, Cambridge, UK. ⁵Max-Planck-Institut für Informatik, Saarbrücken, Germany. ⁶DNAnexus, Mountain View, California, USA. Correspondence should be addressed to E.G. (eg10@sanger.ac.uk) or R.D. (rd109@cam.ac.uk).

Received 1 December 2017; accepted 23 July 2018; published online 20 August 2018; doi:10.1038/nbt.4227

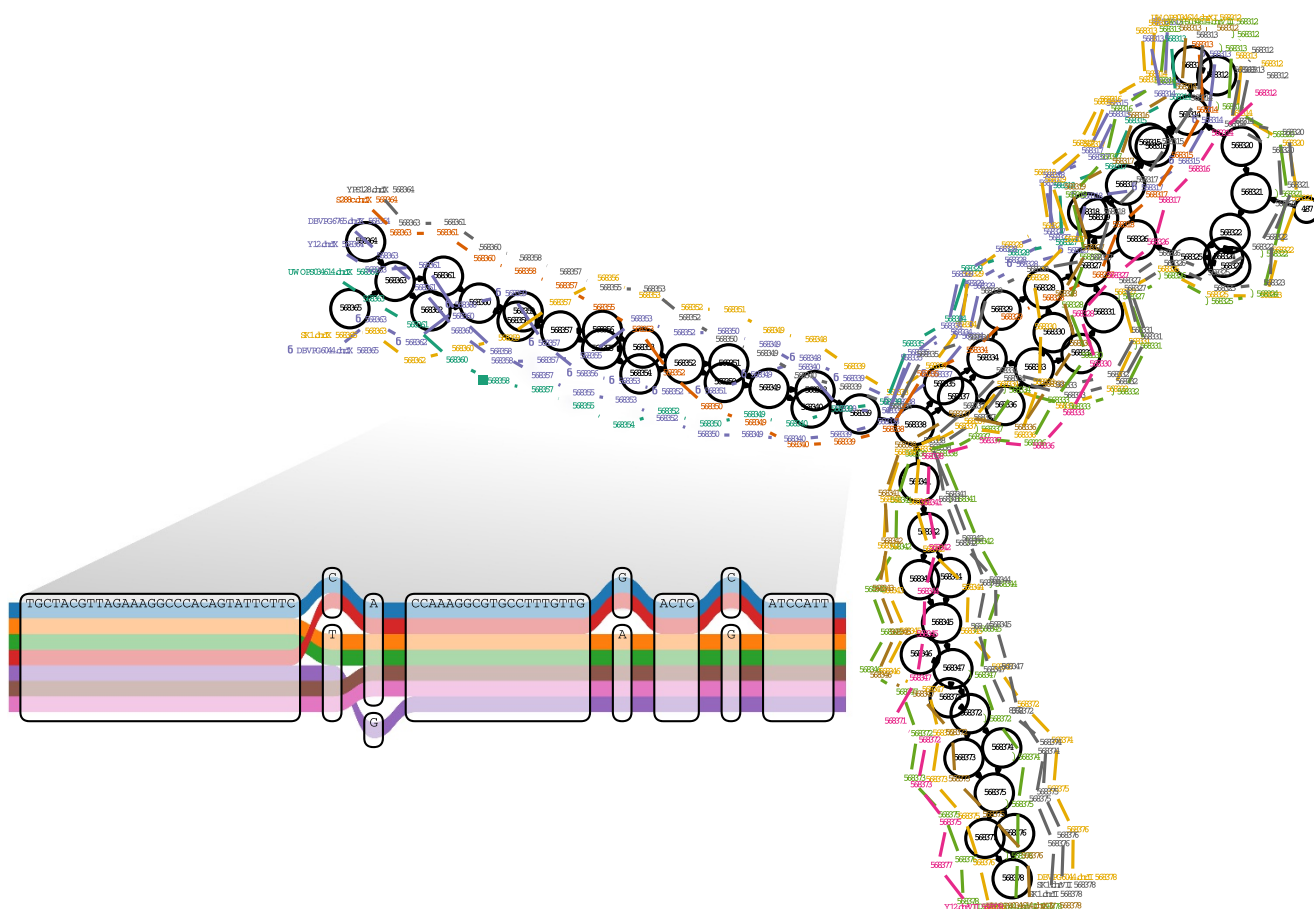


Figure 1 A region of a yeast genome variation graph. This displays the start of the subtelomeric region on the left arm of chromosome 9 in a multiple alignment of the strains sequenced by Yue *et al.*²², built using vg from a full-genome multiple alignment generated with the Cactus alignment package⁷. The inset shows a subregion of the alignment at single-base level. The colored paths correspond to separate contiguous chromosomal segments of these strains. This illustrates the ability of vg to represent paths corresponding to both colinear (inset) and structurally rearranged (main figure) regions of genomic variation.

reference mapping as the current standard, which we take to be bwa mem¹⁸ with default parameters. We show that vg does this, and then that vg maps more informatively around divergent sites.

The final phase of the 1000 Genomes Project (1000GP) produced a data set of ~80 million variants in 2,504 humans¹⁶. We made a series of vg graphs containing all variants or those with minor allele frequency thresholds at 0.1%, 1%, or 10%, as well as a graph corresponding to the standard GRCh37 linear reference sequence without any variation. The full vg graph uses 3.92 GB when serialized to disk, and contains 3.181 Gbp of sequence, which is exactly equivalent to the length of the input reference plus the length of the novel alleles in the VCF file. Complete file sizes including indices range from 25 GB to 63 GB, with details including build and mapping times given in **Supplementary Table 1**.

We next aligned ten million 150-bp paired-end reads simulated with errors from the parentally phased haplotypes of an Ashkenazi Jewish male NA24385, sequenced by the Genome in a Bottle (GIAB) Consortium¹⁹ and not included in the 1000GP sample set, to each of these graphs as well as to the linear reference using bwa mem. **Figure 2a** shows the accuracy of these alignments compared with bwa mem for the 1% allele frequency threshold graph, in terms of receiver operating characteristic (ROC) curves. Comparable plots for other data are given in (**Supplementary Fig. 1**).

Reads that come from parts of the sequence without differences from the reference (middle panel of **Fig. 2a**) mapped slightly better to the reference sequence (green) than to the 1000GP graph (red), which we attribute to a combination of the increase in options for alternative places to map reads provided by the variation graph, and the fact that we needed to prune some search index k-mers in the most complex regions of the graph. As expected, this difference increased as the allele frequency threshold was lowered and more variants were included in the graph (**Supplementary Fig. 1**).

For reads that were simulated from segments containing non-reference alleles (~10% of reads), which are the reads relevant to variant calling, vg mapping to the 1000GP graph (red) gave better performance than either vg (green) or bwa mem (blue) mapping to the linear reference (**Fig. 2b**), because many variants present in NA24385 are already represented in the 1000GP graph. This is particularly clear for single-end mapping, since many paired-end reads are rescued by the mate read mapping. Overall, vg performed at least as well as bwa mem, even on reference-derived reads, and substantially better on reads containing non-reference variants.

We also mapped a real human genome read set with ~50× coverage of Illumina 150-bp paired-end reads from the NA24385 sample to the 1000GP graph. vg produced mappings for 98.7% of the reads, 88.7% with reported mapping quality score 30 on the Phred scale, and 76.8%

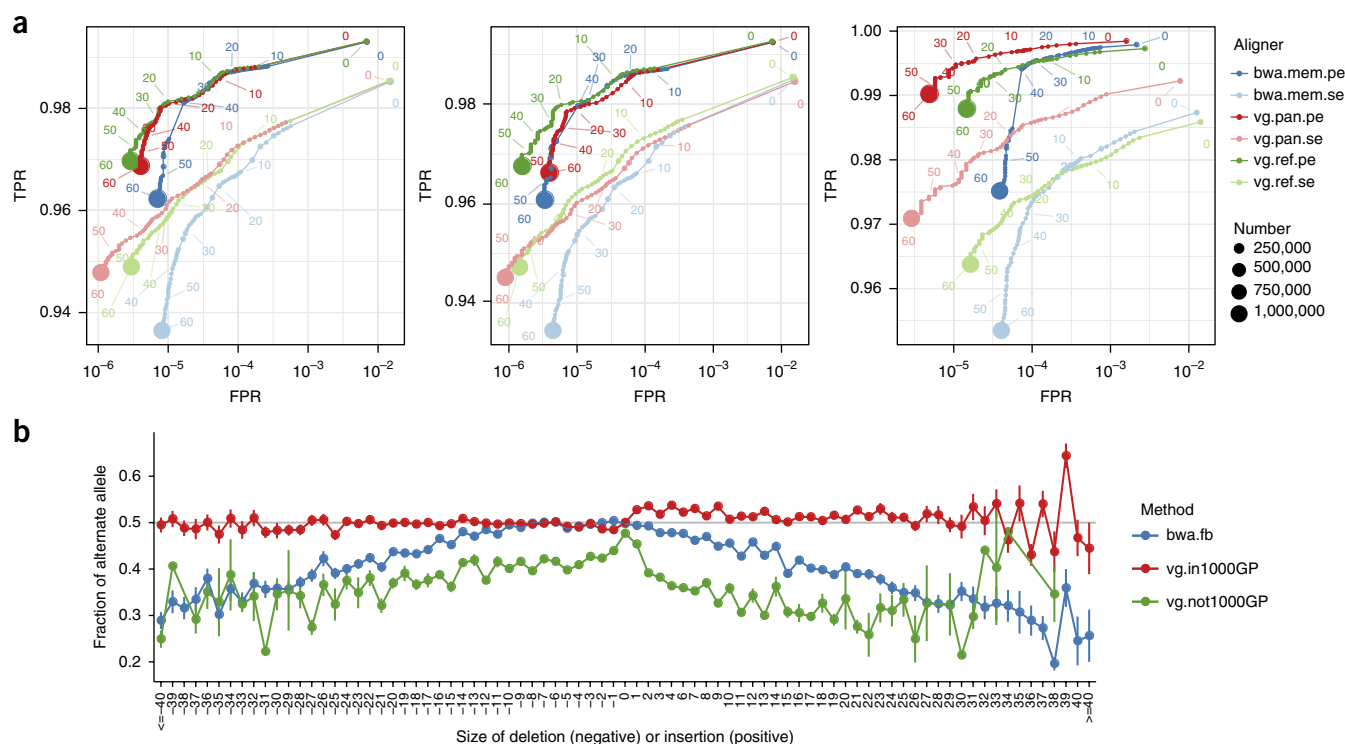


Figure 2 Mapping accuracy for vg against the human genome. **(a)** ROC curves parameterized by mapping quality for 10M read pairs simulated from NA24385 as mapped by bwa mem, vg with the 1000GP 1% allele frequency threshold pangenome reference, and vg with a linear reference, using single-end (se) or paired-end (pe) mapping. Left: all reads, middle: reads simulated from segments matching the linear reference, Right: reads simulated from segments different from the linear reference. **(b)** The mean alternate allele fraction at heterozygous variants previously called in NA24385 as a function of deletion or insertion size (SNPs at 0). Error bars are ± 1 s.e.m.

with perfect, full-length sequence identity to the reported path on the graph. For comparison, we also used vg to map these reads to the linear reference. Similar proportions of reads mapped (98.7%) and with reported quality score 30 (88.8%), but considerably fewer with perfect identity (67.6%). Markedly different mappings were found for 1.0% of reads (0.9% mapping to widely separated positions on the two graphs, and 0.1% mapping to one graph but not the other). The reads mapping to widely separated positions were strongly enriched for repetitive DNA. For example, the linear reference mappings for 27.5% of these read pairs overlapped various types of satellite DNA identified by RepeatMasker, compared to 3.0% of all read pairs.

To illustrate the consequences of mapping to a reference graph rather than a linear reference, we stratified the sites independently called as heterozygous in NA24385 by deletion or insertion length (0 for single-nucleotide variants) and by whether the site was present in 1000GP, and measured the fraction of reads mapped to the alternate allele for each category. The results show that mapping with vg to the population graph when the variant was present in 1000GP (95.4% of sites) gave nearly balanced coverage of alternate and reference alleles independent of variant size, whereas mapping to the linear reference either with vg or bwa mem led to a progressively increasing bias with increasing deletion and (especially) insertion length (Fig. 2b), so that for insertions around 30 bp, a majority of insertions containing reads were missing (there were over twice as many reference reads as alternate reads).

This removal of bias is important when mapping functional genomics data such as ChIP-seq data, where allele-specific expression analysis can reveal genetic variation that affects function but is confounded by reference mapping bias²⁰, especially given that read lengths are

typically shorter for these experiments. We compared mapping with bwa and vg for data set ENCFF000ATK from the ENCODE project²¹, which contains 14.9 million 51-bp ChIP-seq reads for the H3K4me1 histone methylation mark from the NA12878 cell line. When mapping with bwa the ratio of reference to alternate allele matches at heterozygous sites was 1.20, whereas with vg to the 1000GP graph the ratio was 1.01, effectively eliminating reference bias.

We also explored integration of vg with the recently published GraphTyper¹² method, which calls genotypes by remapping reads to a local, partially ordered variation graph built from a VCF file, relying on initial global assignment to a region of the genome by mapping with bwa to a linear reference. Therefore, although GraphTyper also scales to the whole human genome because it is essentially a local method, its functionality is complementary to that of vg, which maps to a global variation graph and does not directly call genotypes. In experiments where we used vg rather than bwa as the primary mapper for GraphTyper, true positives increased marginally (0.02% for single-nucleotide polymorphisms (SNPs) and 0.06% for indels) while false positives increased for SNPs by 0.15% and decreased for indels by 0.03%. We note, however, that GraphTyper was developed by its authors for bwa mapping.

The graphs that we have used so far were constructed from variation data obtained from mapping to a linear reference, and so are directed acyclic graphs. We next demonstrate the ability of vg to work with arbitrary graphs that include duplications, inversions, and translocations, by showing its use with multiple yeast strains independently assembled *de novo* using long-read data²². These assemblies manifest large-scale structural variation and novel sequence not detected in

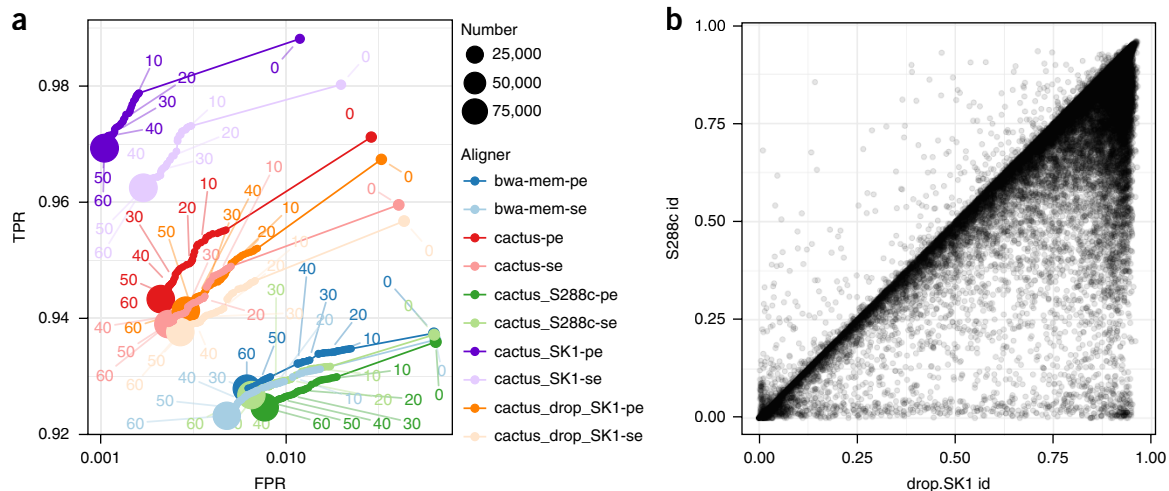


Figure 3 Mapping short and long reads with *vg* to yeast genome references. (a) ROC curves obtained by mapping 100,000 simulated SK1 yeast strain 150-bp paired-end reads against a variety of references described in the text. (b) A density plot of identity fraction when mapping 43,337 Pacific Biosciences long reads from the SK1 strain to the drop.SK1 reference or the S288c reference.

reference-based sequencing, including extensive rearrangement and reordering in subtelomeric regions²² (Fig. 1).

We compared four *vg* graphs: a linear reference graph from the standard S288c strain, a linear reference from the SK1 strain, a pangenome graph of all seven strains, and a “drop SK1” variation graph in which all sequence private to the strain SK1 was removed from the pangenome graph. The multiple genome graphs were constructed with the Cactus progressive aligner⁶, which generates graphs that typically contain cycles and are not partially ordered.

Similarly to the human experiments, we simulated 100,000 150-bp paired-end reads from the SK1 reference, modeling sequencing errors, and mapped them to the four references (ROC curves, Fig. 3a). Not surprisingly, the best performance was obtained by mapping to a linear reference of the SK1 strain from which the data were simulated, with substantially higher sensitivity and specificity compared to mapping to the standard linear reference from the strain S288c with either *vg* or *bwa mem*. Mapping to the variation graphs gave intermediate performance, with >1% more sensitivity and lower false-positive rates than mapping to the standard reference. There was surprisingly little difference between mapping to graphs with and without the SK1 private variation, probably because much of what is novel in SK1 compared to the reference is also seen in other strains. We saw lower sensitivity compared to mapping just to the SK1 sequence, likely because of suppression of GCSA2 index *k*-mers in complex or duplicated regions. In Figure 3b we show the benefit of aligning long reads to a pangenome graph compared to the S288c reference, using a set of 43,337 Pacific Biosciences SK1 reads (mean length 4.7 kb) from reference²².

Finally, to further demonstrate the ability of *vg* to map to arbitrary sequence graphs, we constructed a *vg* graph from a metagenomic assembly of a polar freshwater viral DNA community²³ that was constructed with the *minia3* assembler²⁴. We then aligned a held-out subset of 100,000 reads to this assembly graph using *vg*, and to the linear contigs using *bwa*. Although both methods mapped ~96% of the reads, *vg* had an average identity score of 95% compared to 87% for *bwa*, reflecting that the *bwa* alignments in many cases are not full length (Supplementary Fig. 2).

In conclusion, *vg* implements a suite of tools for genomic sequence data analysis using general variation graph references. Using the *vg*

toolkit, we can construct or import a graph, modify it, visualize it, and use it as a reference. *vg* can accurately map new sequence reads to the reference using succinct indexes of the graph and its sequence space, and can describe variation between a new sample and an arbitrary reference embedded as a path in the graph. Elsewhere¹⁷, we discuss the use of *vg* to map read sets and call variants against a number of alternative human reference graphs built from multiple regions of the human genome with different properties.

There are many areas for potential future development and application of *vg*. These include further improvements in the mapping and variant-calling algorithms, and using long-range statistical haplotype structure information, stored in a graph extension of the positional Burrows–Wheeler transform (PBWT) haplotype compression and search data structure²⁵, as proposed by Novak²⁶. Beyond variant calling, the ability to map in an unbiased way to both reference and alternate alleles is potentially important when quantitating allele-specific protein binding, as shown with ChIP-seq data above, or allele-specific expression²⁷. We note that graphs can also naturally represent the relationships between transcribed, spliced, and edited RNA sequences and the genome from which they are transcribed, so the *vg* software can potentially be used for splice-aware RNA-seq mapping²⁸.

We believe that genome variation graphs will underpin a new paradigm for genome sequence data analysis¹. They support the representation of structural variation using the same components (edges, nodes, and paths) that are used to represent single-base changes. For human, they allow more accurate and complete read mapping (Fig. 2). The benefits will only be greater for other organisms with higher levels of genetic variation, or for which uncertainties remain in the reference assembly. For the biological research community to exploit these advantages, it needs software for variation graphs that scales to the genomes of humans and other complex organisms. *vg* is a robust and openly available platform to fulfill this need.

METHODS

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

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

ACKNOWLEDGMENTS

E.G., J.S., and R.D. were funded by the Wellcome Trust (grants 206194 and 207492). E.T.D. was funded by an NIH Cambridge Trust studentship, and W.J. by a Wellcome Trust MGM studentship (109083/Z/15/Z). A.M.N., G.H., J.M.E., and B.P. were supported by the National Institutes of Health (5U41HG007234), the W.M. Keck Foundation (DT06172015) and the Simons Foundation (SFLIFE# 35190). We thank members of the GA4GH Reference Variation Working Group for support, ideas, and comments, and Hannes Eggertsson for assistance in the integration with GraphTyper.

AUTHOR CONTRIBUTIONS

E.G. conceived and led the development of vg. J.S. developed the GCSA2 index, A.M.N., G.H., J.M.E., and E.T.D. contributed to the software, E.G., W.J., S.G., C.M., M.F.L., and R.D. contributed results and data analysis, R.D. and B.P. oversaw the project, and all contributed to the manuscript.

COMPETING INTERESTS

M.L. is an employee of, and E.G. consults for, DNAnexus Inc. R.D. holds shares in and consults for Congenica Ltd. and Dovetail Inc.

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

- Paten, B., Novak, A.M., Eizenga, J.M. & Garrison, E. Genome graphs and the evolution of genome inference. *Genome Res.* **27**, 665–676 (2017).
- Dilthey, A., Cox, C., Iqbal, Z., Nelson, M.R. & McVean, G. Improved genome inference in the MHC using a population reference graph. *Nat. Genet.* **47**, 682–688 (2015).
- Eggertsson, H.P. *et al.* GraphTyper enables population-scale genotyping using pangenome graphs. *Nat. Genet.* **49**, 1654–1660 (2017).
- Rakocevic, G. *et al.* Fast and accurate genomic analyses using genome graphs. Preprint @bioRxiv <https://doi.org/10.1101/194530> (2017).
- Siren, J. Indexing variation graphs. *Proc. 19th Workshop on Algorithm Engineering and Experiments (ALENEX)* (Society for Industrial and Applied Mathematics, 2017).
- Delcher, A.L. *et al.* Alignment of whole genomes. *Nucleic Acids Res.* **27**, 2369–2376 (1999).
- Paten, B. *et al.* Cactus: Algorithms for genome multiple sequence alignment. *Genome Res.* **21**, 1512–1528 (2011).
- Li, R. *et al.* Building the sequence map of the human pan-genome. *Nat. Biotechnol.* **28**, 57–63 (2010).
- Yuan, S. & Qin, Z. Read-mapping using personalized diploid reference genome for RNA sequencing data reduced bias for detecting allele specific expression. *IEEE International Conference on Bioinformatics and Biomedicine Workshops (IEEE, 2012).*
- Lee, C., Grasso, C. & Sharlow, M.F. Multiple sequence alignment using partial order graphs. *Bioinformatics* **18**, 452–464 (2002).
- Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- Pevzner, P.A., Tang, H. & Waterman, M.S. An Eulerian path approach to DNA fragment assembly. *Proc. Natl. Acad. Sci. USA* **98**, 9748–9753 (2001).
- Myers, E.W. The fragment assembly string graph. *Bioinformatics* **21** (Suppl. 2), ii79–ii85 (2005).
- Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. Preprint @ <https://arxiv.org/abs/1207.3907> (2012).
- DePristo, M.A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- Novak, A.M. *et al.* Genome graphs. Preprint @bioRxiv <https://doi.org/10.1101/101378> (2017).
- Li, H. Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. Preprint @ <https://arxiv.org/abs/1303.3997> (2013).
- Zook, J.M. *et al.* Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
- McDaniell, R. *et al.* Heritable individual-specific and allele-specific chromatin signatures in humans. *Science* **328**, 235–239 (2010).
- ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- Yue, J.-X. *et al.* Contrasting evolutionary genome dynamics between domesticated and wild yeasts. *Nat. Genet.* **49**, 913–924 (2017).
- Aguirre de Cárcer, D., López-Bueno, A., Pearce, D.A. & Alcamí, A. Biodiversity and distribution of polar freshwater DNA viruses. *Sci. Adv.* **1**, e1400127 (2015).
- Chikhi, R. & Rizk, G. Space-efficient and exact de Bruijn graph representation based on a Bloom filter. *Algorithms Mol. Biol.* **8**, 22 (2013).
- Durbin, R. Efficient haplotype matching and storage using the positional Burrows-Wheeler transform (PBWT). *Bioinformatics* **30**, 1266–1272 (2014).
- Novak, A.M., Garrison, E. & Paten, B. in *Algorithms in Bioinformatics* (eds. Firth, M. & Pedersen, C.N.) 246–256 (Springer, Heidelberg, 2016).
- Ge, B. *et al.* Global patterns of cis variation in human cells revealed by high-density allelic expression analysis. *Nat. Genet.* **41**, 1216–1222 (2009).
- Beretta, S. *et al.* in *Algorithms for Computational Biology (AlCoB) 2017*, (eds. Figueiredo, D., Martn-Vide, C., Pratas, D. & Vega-Rodriguez, M.) 49–61 *Lecture Notes in Computer Science* 10252 (Springer, Champaign-Urbana, 2017).

ONLINE METHODS

Model. We define a variation graph to be a graph with embedded paths $G = (N, E, P)$ comprising a set of nodes $N = n_1 \dots n_M$, a set of edges $E = e_1 \dots e_L$, and a set of paths $P = p_1 \dots p_Q$, each of which describes the embedding of a sequence into the graph.

Each node n_i represents a sequence $\text{seq}(n_i)$ that is built from an alphabet $A = \{A, C, G, T\}$. Nodes may be traversed in either the forward or reverse direction, with the sequence being reverse-complemented in the reverse direction. We write n_i^* for the reverse-complement of node n_i , so that $\text{seq}(n_i) = \text{revcomp}(\text{seq}(n_i^*))$; note that $n_i = n_i^{**}$. For convenience, we refer to both n_i and n_i^* as “nodes”. Edges represent adjacencies between the sequences of the nodes they connect. Thus, the graph implicitly encodes longer sequences as the concatenation of node sequences along walks through the graph. Edges can be identified with the ordered pairs of oriented nodes that they link, so we can write $e_{ij} = (n_i, n_j)$. Edges also can be traversed in either the forward or the reverse direction, with the reverse traversal defined as $e_{ij}^* = (n_j^*, n_i^*)$. Note that graphs in vg can contain ordinary cycles (in which n_i is reachable from n_i), reversing cycles (in which both n_i is reachable from n_i^*), and non-cyclic instances of reversal (in which both n_i and n_i^* are reachable from n_j). We implement paths as an edit string with respect to the concatenation of node sequences along a directed walk through the graph. We do not require the alignment described by the edit string to start at the beginning of the sequence of the initial node, nor to terminate at the end of the sequence of the terminal node.

Implementation. The vg implementation is multithreaded and written in C++11, and is available from <https://github.com/vgteam/vg> (version v1.6.0 for code used in the mapping experiments) under the MIT open source software license. It provides both a primary application to support the operations we describe here, and a library libvg, which applications can use to access the data structures, indexes, and low-level operations.

Our core representation of the graph uses Google’s open source protobuf system, which directly supports serialization onto disk for storage. We also provide a protobuf alignment format, GAM (for “Graph Alignment Map”), with analogous functionality to BAM²⁹, but can also export mappings with respect to embedded references in BAM format. To enable read mapping and other random access operations against large sequence graphs we have implemented a succinct representation of a vg variation graph (xg) that is static but very memory and time efficient, using rank/select dictionaries and other data structures from the Succinct Data Structures Library (SDSL)³⁰. Graphs can be imported from and exported to a variety of formats, including the assembly format GFA and the W3C graph exchange format RDF. Further details about the implementation and features are available in the **Supplementary Note** and at the Github website.

Alignment. A key requirement for a reference genome is the ability to efficiently and accurately find an optimal alignment for a new DNA sequence, such as a sequencing read. Analogous to the way that read mappers to linear references work, our approach to this problem is to find seed matches by an indexed search process, cluster them if there are multiple seeds close together, and then perform a local constrained dynamic programming alignment of the read against a region of the graph around each cluster. A brief description of the key steps in this process is given here, with further details in the **Supplementary Note**.

The GCSA2 library⁴ that vg uses for seeding can perform linear time exact match queries independent of the graph size to find super-maximal exact match (SMEM) seeds, subject to a maximum query length, in time comparable to the corresponding operations in bwa mem. SMEMs are exact matches between a query substring and a reference substring that cannot be extended in either direction, and for which there is no extension of the query substring that matches elsewhere in the graph.

After obtaining SMEMs for a query sequence using GCSA2, we cluster them using a global approximate distance metric and distance estimates provided by any nearby paths. For paired reads, we cluster all the SMEMs for both reads in the pair to preferentially support mappings where the SMEMs match a fragment model that we establish online during the alignment of the read set.

We next chain the SMEMs within each cluster by selecting the maximum likelihood path through a Markov model that rewards long SMEMs and short

colinear gaps between SMEMs. In many cases there is just one SMEM in the sequence, but there are complex cases where the best SMEM sequence is not correct, and to catch these we recursively mask out the SMEMs in paths found so far and re-run the algorithm to obtain additional disjoint SMEM sequences if available.

For each consistent sequence of SMEMs, we then obtain the subgraph containing the cluster. To avoid the complications introduced by cycles and inversions³¹, we transform the local graph region into a directed acyclic graph (DAG) while maintaining an embedding in the original, potentially cyclic bidirected graph (**Supplementary Figs. 3 and 4**). We can then perform partial order alignment to the DAG⁹, using banded dynamic programming and an extension of Farrar’s SIMD-accelerated striped Smith–Waterman algorithm³².

When mapping long sequences, we split them into overlapping “chunks” (default 256 bp with 32-bp overlap), map those as above, then chain them using the same colinear Markov model method as described for SMEMs within a chunk. This scales effectively linearly in sequence length up to multiple megabases.

The vg alignment tool also uses base qualities in alignment scores and calculates adjusted mapping quality scores. Base qualities are probabilistic estimates of the confidence of each base call in a read provided by the sequencing technology. vg combines these with a probabilistic interpretation of alignment³³ to adjust the scoring function for alignments, which has previously been shown to improve variant calling accuracy³⁴. Mapping qualities³⁵ are a probability-based measure of the confidence in the localization of a read on the reference that is important for variant calling and other downstream analyses. vg computes mapping qualities by comparing the scores of optimal and suboptimal alignments under the probabilistic alignment model, in a similar fashion to bwa mem.

Graph editing and construction. We can build a graph either by direct construction from external graphs such as from *de novo* assemblies, or by a series of editing operations applied to simple starting graphs such as standard linear reference genomes. To support editing of existing graphs, vg supports operations that can split a node where sequences diverge and insert additional edges and nodes. While doing this it keeps track of the relationship to the previous graph in a translation object, which supports projection of coordinates from one version of the graph to another.

We make use of the editing operations to construct graphs from Variant Call Format (VCF) files¹⁰ as produced by population sequencing projects such as the 1000 Genomes Project¹⁶, inserting a cluster of nodes and edges into a linear reference for each overlapping subset of VCF records. Edit operations also allow progressive construction of a vg graph from a set of sequences by repeated alignment and editing, so that all the initial sequences are embedded in the graph as paths. Last but not least, edit operations allow new variants to be added to an existing vg reference graph to support use cases such as incorporating novel variants from new individuals mapped and called against the graph, while retaining a coordinate mapping to the existing reference. These actions are also invertible, in that vg can generate VCF to describe the graph as a set of variants, using an arbitrarily chosen embedded path as a reference.

Experiments. Experiments were carried out on a dedicated compute node with 256 GB of RAM and two 2.4 GHz AMD Opteron 6378 processors with a total of 32 CPU cores. Mapping comparisons were to bwa version 0.7.15-r1142.

GraphTyper comparisons. To test how vg map complements genotyping in GraphTyper, we mapped reads from the Genome In A Bottle (GIAB) Ashkenazi Jewish Trio benchmark sample HG002 readset, and analyzed variant calling performance on chromosome 21 against the HG002_GRCh37_GIAB_high-conf_CG-III-FB-III-GATKHC-Ion-10X-SOLID_CHROM1-22_v3.3.2_high-conf_triophased.vcf.gz calls using Illumina’s Haplotype Comparison Toolset available from <https://github.com/Illumina/hap.py>. Bwa mem mappings were against the GRCh37d5 reference, and vg mappings against the 1000GP graph then projected onto the GRCh37d5 reference, which is embedded in the 1000GP graph. GraphTyper version 1.3 was run using the dbSNP “common variant” chromosome 21 VCF from NCBI (ftp://ftp.ncbi.nlm.nih.gov/snp/organisms/human_9606_b150_GRCh37p13/VCF/common_all_20170710.vcf.gz). Code used for the analysis is available on request.

Data availability. No new data were collected for this study. The human HG002 data used for **Fig. 2b** are available from ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_son/latest/calls and <http://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP047086> (reads). The yeast whole genome assemblies for **Figures 1** and **3** are available from <http://www.ebi.ac.uk/ena/data/view/PRJEB7245>, the ChIP-seq data set from <https://www.encodeproject.org/files/ENCFF000ATK/>, and the viral metagenome data from <https://www.ebi.ac.uk/ena/data/view/ERS396648>.

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

Code Availability. vg is available at <https://github.com/vgteam/vg> under the MIT open source software license.

29. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
30. Gog, S., Beller, T., Moat, A. & Petri, M. in *International Symposium on Experimental Algorithms* 326–337 (Springer, 2014).
31. Myers, E.W. & Miller, W. Approximate matching of regular expressions. *Bull. Math. Biol.* **51**, 5–37 (1989).
32. Farrar, M. Striped Smith-Waterman speeds database searches six times over other SIMD implementations. *Bioinformatics* **23**, 156–161 (2007).
33. Durbin, R., Eddy, S.R., Krogh, A. & Mitchison, G. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids* (Cambridge University Press, 1998).
34. Hamada, M., Wijaya, E., Frith, M.C. & Asai, K. Probabilistic alignments with quality scores: an application to short-read mapping toward accurate SNP/indel detection. *Bioinformatics* **27**, 3085–3092 (2011).
35. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).

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

There is no new data. We are only using previously reported and publicly available data so this is not relevant.

Data analysis

The vg software is available at <https://github.com/vgteam/vg>. Other software used is previously published and referenced.

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

No new data were collected for this study. The human HG002 data used for figure 2(c) are available from ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_son/latest (calls) and <http://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP047086> (reads). The yeast whole genome assemblies for figures 1 and 3 are available from <http://www.ebi.ac.uk/ena/data/view/PRJEB7245>, the ChIP-seq data set from <https://www.encodeproject.org/files/>

ENCFF000ATK/, the viral metagenome data from <https://www.ebi.ac.uk/ena/data/view/ERS396648> and the NCYC yeast Illumina data are at <http://opendata.ifr.ac.uk/NCYC/>, strains NCYC78, 84, 88, 92, 93, 97, 1006, 1026, 1187, 1228, 1245 and 1681.

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

Life sciences study design

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

Sample size	We simulated 10 million paired reads for the human mapping experiments, and 100,000 for the yeast experiments, and held out 100,000 for the viral metagenome experiments. These at least 10-fold larger than the inverse of the false positive rates we report (1e-6 for human and 1e-3 for yeast), and sufficient to estimate true positive rates to a precision of 0.1% as reported.
Data exclusions	No data were excluded.
Replication	We used held out data, or data simulated from samples not included in building the reference mapped to.
Randomization	Randomization was not relevant to this study.
Blinding	There was no blinding. Experiments were computational and all results reported, so there was no human component to the numerical results reported, so no requirement for blinding.

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 checked="" type="checkbox"/>	<input 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