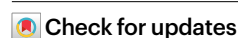


# Haplotype-aware pantranscriptome analyses using spliced pangenome graphs

Received: 18 June 2021

Accepted: 28 November 2022

Published online: 16 January 2023



Jonas A. Sibbesen<sup>1,3</sup>, Jordan M. Eizenga<sup>1,3</sup>, Adam M. Novak<sup>1</sup>, Jouni Sirén<sup>1</sup>, Xian Chang<sup>1</sup>, Erik Garrison<sup>2</sup> & Benedict Paten<sup>1</sup>✉

Pangenomics is emerging as a powerful computational paradigm in bioinformatics. This field uses population-level genome reference structures, typically consisting of a sequence graph, to mitigate reference bias and facilitate analyses that were challenging with previous reference-based methods. In this work, we extend these methods into transcriptomics to analyze sequencing data using the pantranscriptome: a population-level transcriptomic reference. Our toolchain, which consists of additions to the VG toolkit and a standalone tool, RPVG, can construct spliced pangenome graphs, map RNA sequencing data to these graphs, and perform haplotype-aware expression quantification of transcripts in a pantranscriptome. We show that this workflow improves accuracy over state-of-the-art RNA sequencing mapping methods, and that it can efficiently quantify haplotype-specific transcript expression without needing to characterize the haplotypes of a sample beforehand.

Transcriptome profiling by RNA sequencing (RNA-seq) has matured into a standard and essential tool for investigating cellular state. Bioinformatics workflows for processing RNA-seq data generally begin by comparing reads to a reference genome or reference transcriptome<sup>1–4</sup>. This is an expedient method that makes it practical to analyze the large volume of data produced by high-throughput sequencing.

Reference-based methods also have costs. When the genome of a sample differs from the reference, bioinformatics tools must account for the resulting mismatches between the sequencing data and the reference. This results in reduced ability to correctly identify reads with their transcript-of-origin, with larger genomic variation leading to a greater reduction in accuracy. This problem is known as reference bias<sup>5</sup>.

Computational pangenomics has emerged as a powerful methodology for mitigating reference bias. Pangenomics approaches lean heavily on abundant, publicly available data about common genomic variation for certain species (notably including humans). These methods incorporate population variation into the reference itself, usually in the form of a sequence graph<sup>6</sup>. Mapping tools for pangenomic references have demonstrated reduced reference bias when mapping DNA reads<sup>7,8</sup>. This facilitates downstream tasks that are frustrated by mapping biases, such as structural variant calling<sup>9,10</sup>.

The sequence graph formalism used in pangenomics has an additional attractive feature for RNA-seq data: it can represent splice junctions with little modification. Without this benefit, RNA-seq mappers for conventional references must make use of sometimes elaborate algorithmic heuristics to align over known splice junctions<sup>2</sup>. Alternatively, they can map to only known isoforms, but this technique introduces mapping ambiguity owing to the reuse of exons across isoforms<sup>11</sup>.

The current methodological landscape in pangenomics is ripe to be extended to pantranscriptomics: using populations of reference transcriptomes to inform transcriptomic analyses. There is some precedent in previous transcriptomic methods that have used sequence graphs. AERON<sup>12</sup> uses splicing graphs and GRAPHALIGNER<sup>13</sup> to identify gene fusions. ASGAL<sup>14</sup> uses splicing graphs to identify novel splicing events. Finally, the pangenomic aligner HISAT2<sup>15</sup> is built on the RNA-seq aligner HISAT<sup>16</sup> and retains many of its features for RNA-seq data.

One transcriptomic analysis that is particularly prone to reference bias is allele-specific expression (ASE). ASE seeks to identify differences in gene expression between the two copies of a gene in a diploid organism. These differences are indicative of various biological processes, including *cis*-acting transcriptional regulation, nonsense-mediated decay, and genomic imprinting<sup>17,18</sup>. The differences are identified by

<sup>1</sup>UC Santa Cruz Genomics Institute, Santa Cruz, CA, USA. <sup>2</sup>University of Tennessee Health Science Center, Memphis, TN, USA. <sup>3</sup>These authors contributed equally: Jonas A. Sibbesen and Jordan M. Eizenga. ✉e-mail: [bpaten@ucsc.edu](mailto:bpaten@ucsc.edu)

measuring the ratio between RNA-seq reads containing each allele of a heterozygous variant. However, the reads containing the non-reference allele are systematically less mappable because of reference bias, which can confound signals of ASE<sup>5</sup>. Several approaches have been developed to deal with reference bias for ASE detection. WASP filters reads that show allele-biased mapping before ASE estimation<sup>19</sup>. Others can mitigate bias at the read mapping stage, but require variant calls, often with phasing, for the individual being analyzed<sup>20,21</sup>. The variant information is either incorporated into the mapping algorithm to reduce reference bias or used to create a sample-specific diploid reference to map against.

Pantranscriptomic approaches using existing haplotype panels for inferring haplotype-specific expressions in smaller regions have also been developed specifically for the human leukocyte antigen (HLA) region. ALTHAPALIGNR and HLAPEERS both align reads to a set of HLA haplotypes<sup>22,23</sup>. The alignments are then used to infer haplotype-specific gene or transcript expression.

In this work, we present a bioinformatics toolchain for genome-wide pantranscriptomic analysis, which consists of additions to the VG toolkit and a standalone tool, RPVG. First, VG RNA can combine genomic variation data and transcript annotations to construct a spliced pangenome graph. Next, VG MPMAP can align RNA-seq reads to these graphs with high accuracy. Finally, RPVG can use the alignments from VG MPMAP to quantify haplotype-specific transcript expression. The population variation that is embedded in the pantranscriptome reference makes it possible to do so without first characterizing the sample genome and without restricting focus to single-nucleotide variants (SNVs).

## Results

### Haplotype-aware transcriptome analysis pipeline

In short, our pipeline works as follows. First, we construct a spliced pangenome graph and a pantranscriptome using VG RNA, a tool developed as part of the VG toolkit<sup>7</sup> (Fig. 1a). The pantranscriptome consists of a set of haplotype-specific transcripts (HSTs) and is constructed by projecting (lifting over) the transcripts in a transcript annotation onto a set of known haplotypes. VG RNA uses the graph Burrows–Wheeler transform (GBWT) to efficiently store the HST paths, allowing the pipeline to scale to a pantranscriptome with millions of transcript paths<sup>24</sup>. Next, RNA-seq reads are mapped to the spliced pangenome graph using VG MPMAP, a new splice-aware graph mapper that can align across both annotated and unannotated splice junctions (Fig. 1b). VG MPMAP produces multipath alignments that capture the local uncertainty of an alignment to different paths in the graph (Extended Data Figure 1). Lastly, the expression of the HSTs are inferred from the multipath alignments using RPVG (Fig. 1c). RPVG uses a nested inference scheme that first infers the most probable underlying haplotype pairs and then estimates the HST expression using expectation maximization.

### RNA-seq mapping benchmark

We compared VG MPMAP against three other mappers: STAR<sup>2</sup>, HISAT2<sup>15</sup> and VG MAP<sup>7</sup>. STAR and HISAT2 can both use splicing information to guide mapping. However, of the two, only HISAT2 is able to also utilize genomic variants. VG MAP is not a splice-aware mapper, but it is still able to map to spliced pangenome graphs, which contain both splicing and genomic variation edges.

We used two different references for the comparison: the standard reference genome with added splice junctions (spliced reference) and a spliced pangenome graph containing both splice junctions and variants (spliced pangenome graph). For STAR, only the spliced reference was used. In addition, to assess alignment across unannotated splice junctions, we constructed references with a random 20% of transcripts removed before construction (on the basis of recent estimates of the fraction of novel transcripts in a sample<sup>25</sup>). For all of the tools besides

STAR, this reference included variation (80% spliced graph), whereas the reference for STAR did not (80% spliced reference).

**Simulated sequencing data.** Paired-end reads were simulated from HSTs derived from the GENCODE transcript annotation set<sup>26</sup> and the NA12878 haplotypes from the 1000 Genomes Project (1000GP)<sup>27</sup> using VG SIM. The CEU population was excluded from the spliced pangenome graph, as NA12878 is from that population, and we wanted to estimate performance for an individual who may not be closely related to the 1000GP populations.

Using the set of simulated reads we first compared the overall mapping performance of each method. Figure 2a shows the mapping error (1 – precision) and recall for different mapping quality thresholds.

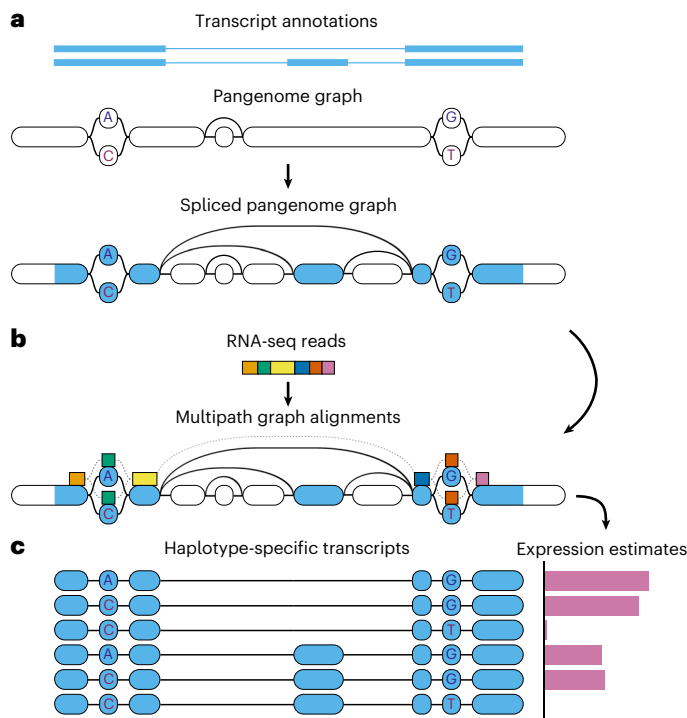
Reads are considered correctly mapped if one of their multi-alignments covers 90% of the true reference sequence alignment. As can be seen in Fig. 2a, VG MPMAP achieves both a low error and high recall, while the other methods either had a high error or low recall. The same pattern is observed among primary alignments, ignoring multi-alignments (Extended Data Figure 2). The results also show that the spliced pangenome graph generally improves mapping performance. In addition, VG MAP and VG MPMAP show substantially better calibration in their estimated mapping qualities, especially among the most confidently mapped reads (Supplementary Figure 1).

On the 80% spliced references, the performance of all of the tools decreases relative to the corresponding reference constructed with the full transcript set. As expected, the performance of VG MAP decreases dramatically, as it can only align over splice junctions represented in the graph. The reduction in the performance of VG MPMAP is larger than for STAR and HISAT2. This reduction is concentrated on reads containing unannotated splice junctions (Supplementary Figure 2), but the performance of VG MPMAP is still competitive with both of the other tools in the aggregate read set.

Using a fixed mapping quality threshold, we evaluated how the methods perform for different edit distances between the simulated reads and the reference. Extended Data Figure 3 shows this analysis for unique (mapping quality of at least 30) and multi-alignments. VG MPMAP achieves a high recall even for reads with an edit distance above 3. The recall of HISAT2, and to a lesser extent STAR, markedly decreases for the same distance.

Next, we evaluated whether using a variant-aware approach reduces reference bias. Figure 2b shows the mean fraction of reads mapped to the alternative allele for different allele lengths. When using the spliced reference genome, all methods exhibit a bias towards the reference allele, with VG MAP and MPMAP showing less bias than the other methods. Using the spliced pangenome graph results in substantially reduced bias for all methods. We also analyzed the mapping error and recall on reads stratified by the number of variants they contain (Extended Data Figure 4). This analysis corroborates the allele bias results; VG MPMAP and VG MAP retain high recall in the presence of variants, whereas the performance of HISAT2 and STAR decreases substantially, especially in the presence of indels.

Previous research has pointed out that allelic bias can also result from differential uniqueness between two alleles<sup>19</sup>. The WASP tool combats this bias by filtering out potentially biased reads. We compared allelic bias between the four mapping tools and a pipeline consisting of WASP and STAR. Using simulated data with no allelic bias, we identified heterozygous variant sites with coverage at least 20 and measured (i) the number of such sites and (ii) the proportion of sites with a statistically significant allele skew (two-sided binomial test,  $\alpha = 0.01$ ) (Extended Data Figure 5). Both HISAT2 and STAR show an increase in falsely significant tests above the nominal false positive rate of 0.01, especially for insertions and deletions. The WASP (STAR) pipeline, VG MAP, and VG MPMAP all show approximately the expected rate of false positives for all variant types. In addition, compared to the



**Fig. 1 | Diagram of haplotype-aware transcriptome analysis pipeline.** The three major steps in the pipeline. **a**, VG RNA adds splice junctions derived from a transcript annotation to a pangenome graph to create a spliced pangenome graph. It simultaneously creates a pantranscriptome composed of a set of haplotype-specific transcripts (HSTs) using a panel of known haplotypes (not shown). **b**, VG MPMAP aligns RNA-seq reads to subgraphs of the spliced pangenome graph represented as a multipath alignment. **c** RPVG uses the alignments from MPMAP to estimate the expression of the HSTs in the pantranscriptome.

WASP (STAR) pipeline, VG MPMAP retains 5,670 more variant sites with coverage at least 20.

The mapping results were corroborated by alternative evaluation methodologies. First, we used an alternate correctness criterion based on aligning within 100 bases of the correct position on the paths in the graph (Supplementary Figure 3), which gave qualitatively similar results. Second, we used RSEM as an alternative read simulator (Supplementary Figure 4). All mapping tools showed similar performance with the alternative simulator except HISAT2, which had relatively higher recall.

The set of simulated reads used for the mapping evaluation presented in Fig. 2a,b was not used to optimize the algorithmic design or parameters of VG MAP and VG MPMAP. Thus, these reads can be considered a ‘test set’. Supplementary Figure 5a–c shows the results on one of the simulated ‘training sets’ that were used to optimize the method. Simulated data from RSEM was also used during the development of VG MPMAP.

**Real sequencing data.** We used RNA-seq reads from the ENCODE project (ENCSR000AED) to benchmark the methods on real data<sup>28,29</sup>. We first looked at the fraction of aligned reads for each method (Fig. 2c). As can be seen in the figure, all methods have comparable overall mapping rates. When only looking at alignments with a mapping quality value of at least 30, both STAR and HISAT2 show noticeably higher rates compared to VG MAP and VG MPMAP. However, it seems that the cost of these higher mapping rates is lower precision (Fig. 2a) and poorly estimated mapping qualities (Supplementary Figure 1).

Ground-truth alignments are not available for real data, so instead we use a proxy that is based on Pacific Biosciences (PacBio) Iso-Seq read

alignments generated by the ENCODE project (ENCSR706ANY), which are from the same cell line. We expect the transcript expression to be similar despite some technical biases owing to the different sequencing protocols, and long reads can be mapped more confidently than short reads. Thus, higher correlation in coverage between the mappings should be indicative of more accurate short read mappings. Figure 2d shows the estimated Pearson’s correlation in the coverage of each exon as a function of mapping quality threshold. Exons were defined by the Iso-Seq alignments. As can be seen, both VG MAP and VG MPMAP achieves higher correlation than STAR and HISAT2, with the spliced pangenome graph resulting in even higher correlation for both (see Supplementary Figure 6 for the full scatter plot).

Finally, we compared the methods’ computational requirements. Figure 2e shows the number of read pairs mapped per second per thread. Conversion from SAM to BAM was included in the HISAT2 time estimate to be more comparable to the output type of the other methods. VG MPMAP is 3.1–4.6 times slower than HISAT2, depending on the graph, but 10.3 times faster than VG MAP on the spliced pangenome graph. VG MPMAP uses slightly more memory than STAR on the spliced reference and somewhat more memory than HISAT2 on the spliced pangenome graph (Fig. 2f).

Results on additional datasets used during the development of VG MPMAP can be seen in Supplementary Figs 5d,e and 7. The same data were also used to optimize the parameters of VG MAP for RNA-seq mapping.

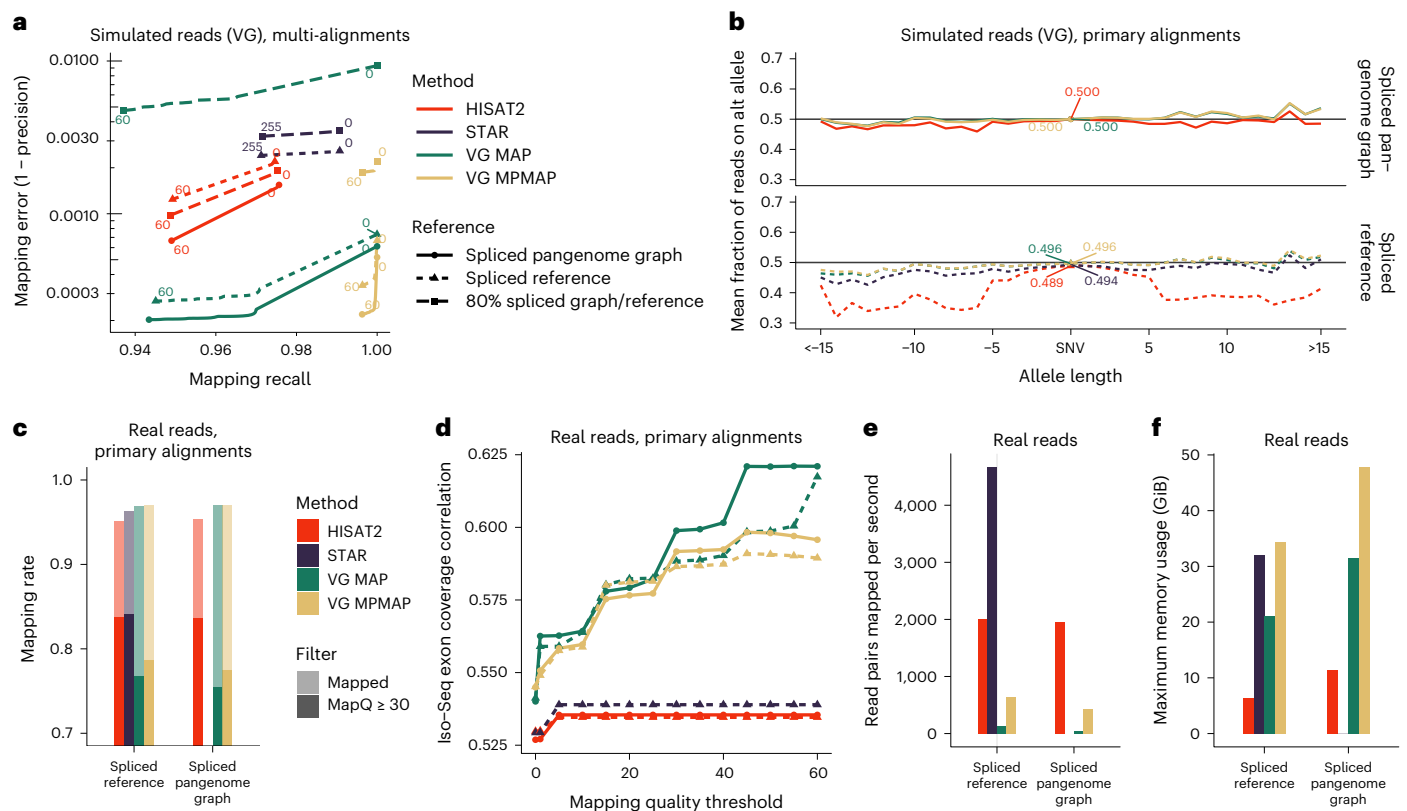
### Haplotype-specific transcript quantification

We compared RPVG to three other transcript quantification methods: KALLISTO<sup>3</sup>, SALMON<sup>4</sup>, and RSEM<sup>1</sup>. We stress that none of these methods were developed to work on pantranscriptomes with millions of HSTs. However, they serve as a point of reference for what accuracy is achievable without new methods development. RPVG’s inference model includes both a diplotype and HST expression, conditioned on the diplotype. However, to facilitate the comparison to other tools, we report here the marginal expression over all HSTs, which is directly comparable to the output of the other tools that lack a diplotype model.

Three different pantranscriptomes were generated for the evaluation using different sets of 1000GP haplotypes (Supplementary Table 3): (i) all European haplotypes excluding the CEU population (“Europe (excl. CEU)”, 2,515,408 HSTs); (ii) all haplotypes excluding the CEU population (“Whole (excl. CEU)”, 11,626,948 HSTs); and (iii) all haplotypes (“Whole”, 11,835,580 HSTs). The CEU population was excluded for the same reason as in the mapping benchmark. In addition, we created a personal-sample-specific transcriptome consisting of NA12878 HSTs (“Personal (NA12878)”, 235,400 HSTs). This transcriptome corresponds to the ideal case where a sample’s haplotypes are known beforehand. HSTs with a haplotype probability below 0.8 were filtered from the RPVG output (Supplementary Figure 8).

**Simulated sequencing data.** We first looked at the ability of the method to accurately predict whether an HST was expressed or not. Figure 3a shows the recall and precision using simulated data. The results were stratified by different expression thresholds up to a value of 10 TPM (transcripts per million). Note that we were not able to run RSEM on the two largest pantranscriptomes used in the analysis. RPVG exhibits a much higher precision than the other tools for all pantranscriptomes. This illustrates the importance of having a diplotyping model when inferring HST expression using a pantranscriptome reference, which is one of the major differences between RPVG and the other methods. Importantly, only a minor difference is observed between the pantranscriptomes without the CEU population (excl. CEU) and the whole pantranscriptome (“Whole”), which contains NA12878. This could be explained by the fact that less than 2% of HSTs are on average unique to a specific sample when compared to all samples in other populations using the 1000GP data (Extended Data Figure 6). This





**Fig. 2 | Mapping benchmark using RNA-seq data from NA12878.** RNA-seq mapping results comparing VG MPMAP and three other methods using simulated and real Illumina data. **a**, Mapping error and recall for different mapping quality thresholds (colored numbers) using simulated data. Reads are considered correctly mapped if one of their multi-alignments covers 90% of the true reference sequence alignment. **b**, Mean fraction of mapped reads supporting the non-reference allele for variants of different lengths in simulated data. Negative

lengths correspond to deletions and positive to insertions. The colored numbers are the mean fraction for SNVs. **c**, Mapping rate using real data. **d**, Pearson's correlation between Illumina and Iso-Seq exon coverage using real data as a function of mapping quality threshold. **e**, Number of read pairs mapped per second per thread using real data on an AWS m5.4xlarge instance. **f**, Maximum memory usage for mapping in gigabytes using real data.

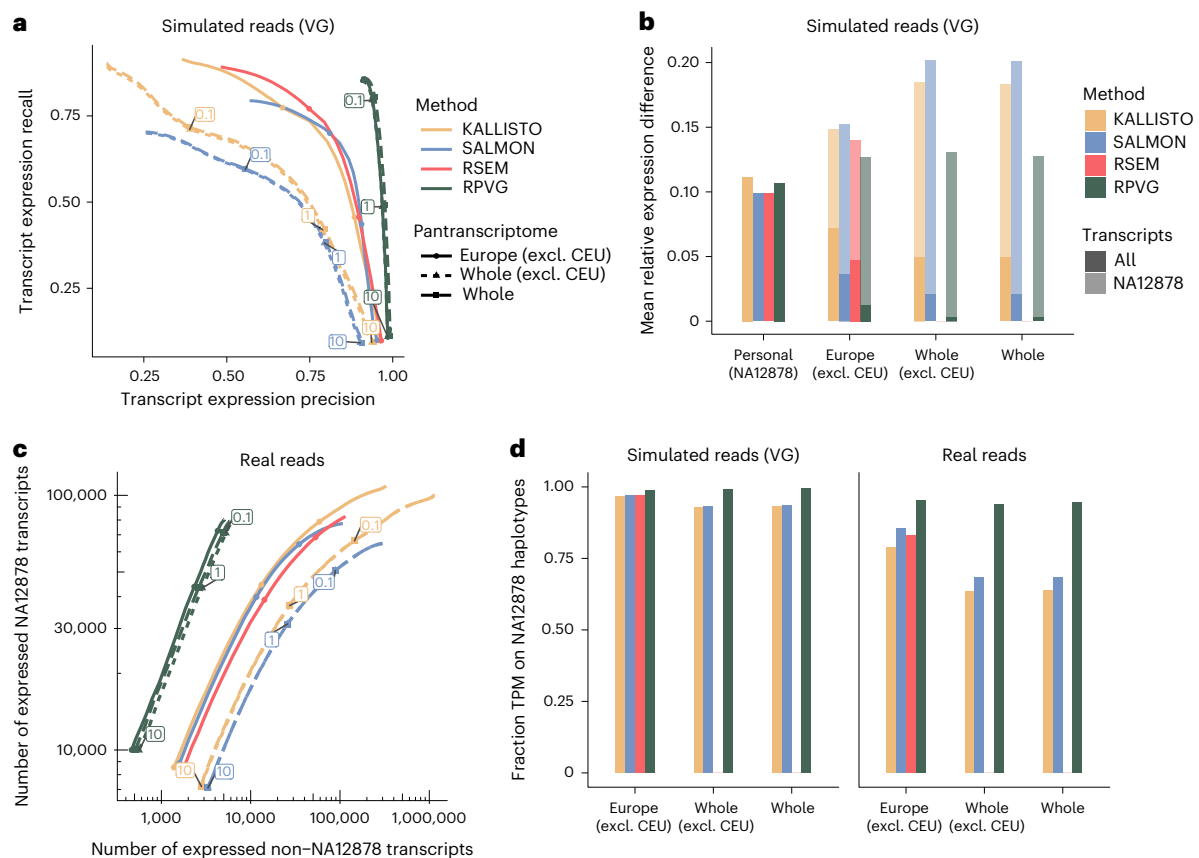
suggests that haplotype panels like the 1000GP are a good alternative when a sample's haplotypes are not available.

We compared how well the different methods could predict the correct expression value. Figure 3b shows the mean absolute relative expression difference (MARD) between the expression values of the simulated reads and the estimated values. On the personal set, RPVG performs comparably to the other methods. However, as the size of the pantranscriptome grows, the MARD on the NA12878 transcript set only increases slightly for RPVG. Supplementary Figure 9 shows the full scatter plots of the simulated and estimated expression values for the NA12878 HSTs. The lower error for all methods when using all HSTs can be explained by the larger number of unexpressed HSTs. Comparing Spearman correlations gives similar conclusions, except that KALLISTO and RSEM perform comparably to RPVG when restricting focus to the haplotypes of NA12878 (Supplementary Figure 10). This suggests that KALLISTO and RSEM accurately ranks the expression of these transcripts but do not accurately estimate the absolute quantity. Using only the HSTs estimated to be expressed by each method, we see similar results for MARD and Spearman correlation (Supplementary Figure 11). However, when looking at reference transcript-level expression estimates by summing over HSTs, the other methods exhibit overall better MARDs (Supplementary Figure 12).

RPVG can optionally use Gibbs sampling to quantify the uncertainty in the expression estimates. To evaluate the accuracy of this procedure, we estimated 90% credible intervals from 1,000 samples per HST (Supplementary Figure 13); 86.4% of the intervals contained the simulated expression value, which is close to the expected proportion.

We also compared the ability of the VG MPMAP–RPVG pipeline to estimate ASE to a pipeline of WASP<sup>19</sup> with STAR<sup>2</sup> alignments. This analysis focused on allele-specific read counts over heterozygous variants. We converted the simulated HST expression values to read counts and defined true positives as variants with significant ASE using a two-sided binomial test with *P* values adjusted using the Benjamini–Hochberg procedure with false discovery rate (FDR)  $\alpha = 0.1$ . We converted the RPVG estimates into read counts similarly and then called ASE with the read counts of both pipelines using the same statistical procedure as with the simulated values (Extended Data Figure 7). The VG MPMAP–RPVG pipeline achieves a markedly higher true positive rate with the same false positive rate as the WASP–STAR pipeline. Moreover, VG MPMAP–RPVG had similar performance for indels, whereas WASP excludes these variants.

**Real sequencing data.** Next, we evaluated the accuracy of the HST expression estimation using real sequencing data from the ENCODE project (ENCSTR000AED)<sup>28,29</sup>. Since we do not know which transcripts are expressed in real data, we focus instead on the haplotype estimation. We can indirectly measure accuracy by asking whether the HSTs that are estimated to be expressed are in fact from NA12878. Figure 3c shows that RPVG predicts markedly fewer HSTs from non-NA12878 haplotypes than the other methods. Also, we see again only a minor difference between pantranscriptomes. Next, we compared the fraction of transcript expression (in TPM) that was attributed to NA12878 haplotypes for simulated (left) and real (right) data (Fig. 3d). RPVG attributes more than 98.8% and 94.0% of the expression to NA12878



**Fig. 3 | HST quantification benchmark using RNA-seq data from NA12878.** HST quantification results comparing RPVG against three other methods using simulated and real Illumina data. It should be noted that the other methods were primarily designed for reference transcript quantification and not millions of HSTs. For details on the pantranscriptomes used see Supplementary Table 3. **a**, Recall and precision of whether a transcript is correctly assigned non-zero expression for different expression value thresholds in transcripts per million (TPM; colored numbers for “Whole (excl. CEU)”) using simulated data. **b**, MARD between simulated and estimated expression (in TPM) for different

pantranscriptomes using simulated data. MARD was calculated using either all HSTs in the pantranscriptome (solid bars) or using only the NA12878 HSTs (shaded bars). “Personal (NA12878)” is a sample-specific transcriptome. **c**, Number of expressed transcripts from NA12878 haplotypes against the number from non-NA12878 haplotypes for different expression value thresholds (colored numbers) using real data. **d**, Fraction of transcript expression (in TPM) assigned to NA12878 haplotypes for different pantranscriptomes using simulated (left) and real (right) data.

haplotypes when using simulated and real data, respectively. Furthermore, the prediction accuracy of RPVG only decreases slightly when the size of the pantranscriptome increases from 2.5 million HSTs in “Europe (excl. CEU)” to 11.6M in “Whole (excl. CEU)”.

To assess the robustness of the VG MPMAP–RPVG pipeline to samples with recently admixed ancestry, we applied it to two samples from a recent study<sup>30</sup>: one of European American ancestry and one of African American ancestry (Extended Data Figure 8). We expect that the African American individual has a more admixed ancestry owing to the greater genomic diversity present in Africa and the history of widespread slave rape by slave owners of European ancestry in the United States<sup>31</sup>. As a proxy for accuracy, we quantify how frequently RPVG can identify two or fewer HSTs for a transcript (if none of the HSTs match the individual, the posterior will tend to diffuse onto multiple similar HSTs). Consistent with expectations from Extended Data Figure 6, we see somewhat lower accuracy on the African American individual, but the difference is small.

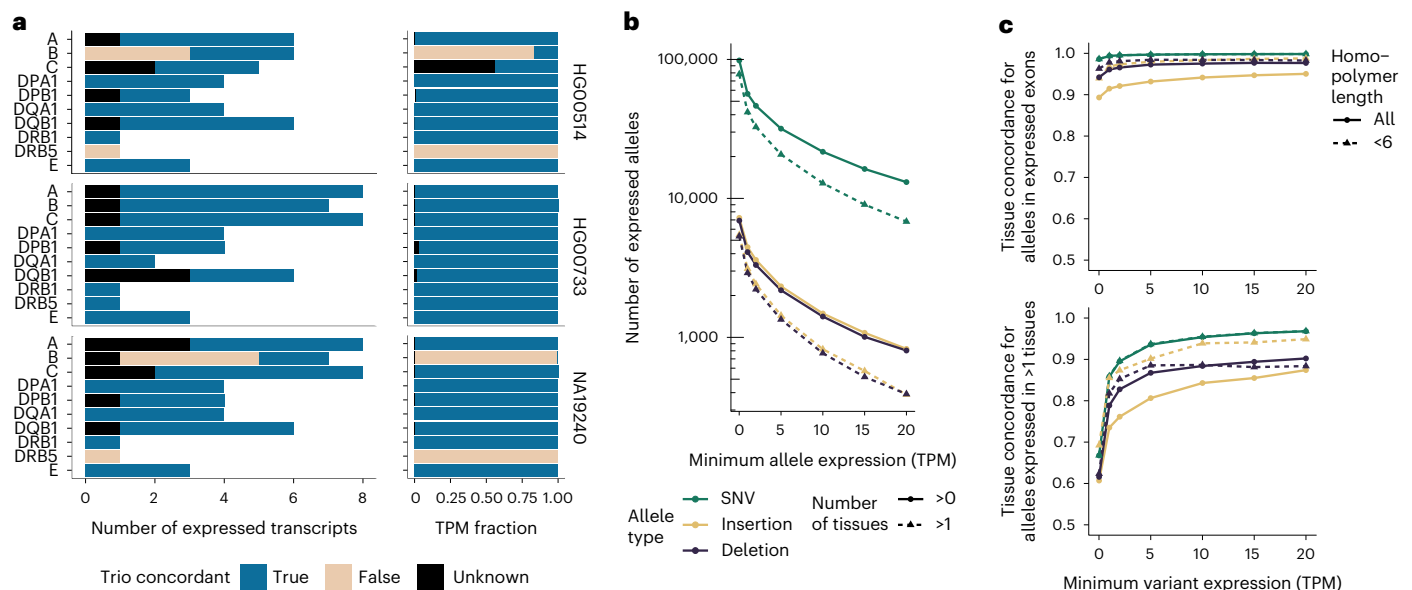
To show the advantage of multipath alignments for inference, we repeated the simulated and real data evaluations using single-path alignments from VG MPMAP (taking the best scoring path in each multipath alignment) and VG MAP (Extended Data Figure 9). For all pantranscriptomes and datasets, RPVG gave the best results using the multipath alignments.

Results on additional simulated and real datasets used when developing RPVG, including selecting its default parameters, can be seen in Supplementary Figs 14 and 15.

### Evaluating HLA typing

We evaluated the ability of the VG MPMAP–RPVG pipeline to infer diplotypes for genes in the highly polymorphic HLA region. To do so, we created two HLA-specific pantranscriptomes using the IPD-IMGT/HLA database<sup>32</sup> (see Supplementary Table 3). We ran the pipeline on RNA-seq data for three parent–child trios from the 1000GP sequenced in the Human Genome Structural Variation Consortium (HGSVC)<sup>33</sup> (see Supplementary Table 4). Figure 4a shows the number of predicted expressed transcripts for each child and the Mendian concordance of the inferred parent and child diplotypes. The same results are also summarized by proportion of inferred expression. With the exception of B and DRB5, almost all of the genes’ expression is assigned to concordant transcripts.

We also ran the pipeline on ten randomly selected CEU samples from Geuvadis<sup>34</sup> for which HLA typing results are available from other studies of genomic sequencing data<sup>35,36</sup>. The results were similar to the trios: A, DQB1, and DRB1 had correct typing in all samples, and B was incorrect in some of the samples (Supplementary Figure 16). However, typing of C was also incorrect in some of the Geuvadis samples.



**Fig. 4 | HLA typing and allele concordance evaluation using RNA-seq data from trios and different tissues.** **a**, Mendelian concordance of HLA typing results using Illumina data from three trios and a pantranscriptome containing ten HLA genes. Results are summarized by the number of transcripts (left) and proportion of expression in TPM (right) predicted to be expressed for each child and gene. The concordance is labeled unknown when a transcript is not expressed in one of the parents. **b**, **c**, Variant genotyping analyses using Illumina data from five tissues from the same individual and a pantranscriptome

containing the 1000GP haplotypes. **b**, Number of variant alleles predicted to be expressed in at least one (solid lines) or two tissues (dashed lines) for different expression thresholds. **c**, Fraction of alleles predicted to be concordant across tissues for alleles in all expressed exons (including unexpressed alleles of expressed variants; top) and alleles expressed in at least two tissues (bottom). The results are shown for different variant expression thresholds and homopolymer lengths. See Extended Data Figure 10 for a graphical description of concordance.

While the results look promising, other HLA typing methods have shown similar or somewhat higher accuracy, depending on the gene, although the small sample size makes it challenging to determine the exact difference<sup>37</sup>. However, one major advantage of the VG MPMAP–RPVG pipeline compared to these methods is that it also provides HST expression estimates in addition to the typing.

### Investigating variant genotyping and effect prediction

To illustrate the ability of the VG MPMAP–RPVG pipeline to genotype variants in a pantranscriptome from RNA-seq data, we ran the pipeline on five different tissue samples from the same individual, sequenced by the ENCODE project<sup>28,29</sup> (see Supplementary Table 4). Figure 4b shows the number of expressed variant alleles for different expression thresholds. As expected, markedly more SNVs are predicted to be expressed than indels. A similar number of insertions and deletions are predicted to be expressed.

For validation, we looked at whether the inferred alleles were concordant across tissues. An allele was considered concordant if it was either consistently expressed or consistently not expressed across all tissues for which the corresponding variant is expressed (see Extended Data Figure 10 for a graphical description). To account for allelic drop-out for lowly expressed exons, we calculated the concordance for different thresholds of total variant expression. Figure 4c shows the results of this analysis for alleles in all expressed exons (including unexpressed alleles of expressed variants; top) and alleles expressed in at least two tissues (bottom). Across all expressed exons, the concordance rate reaches 0.95 for insertions, with higher values for deletions and SNVs. For alleles expressed in at least two tissues, the rates are lower but still over 0.95 for SNVs and 0.85 for indels. After removing variants in homopolymers longer than five bases, the performance on insertions improves substantially, although, surprisingly, the performance for deletions was largely unchanged.

Finally, we investigated the effect of the predicted variants on functional elements using the Ensembl Variant Effect Predictor (VEP)

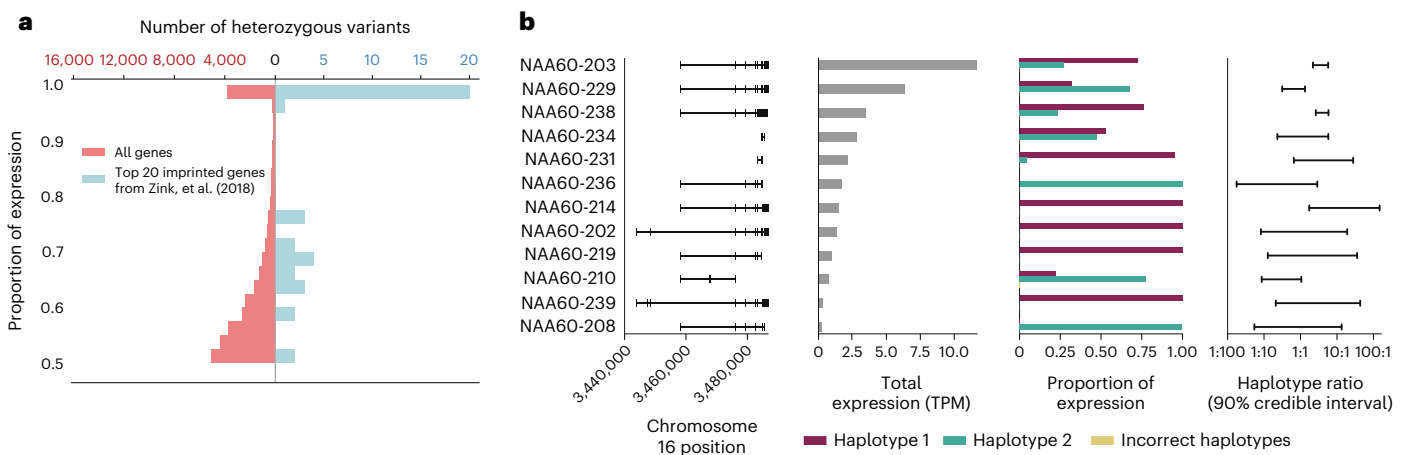
toolset<sup>38</sup> (Supplementary Figure 17). Among variants in exons with TPM of at least five, the number of predicted protein truncating variants (frameshift, splice donor, splice acceptor, and stop gain) is comparable to or lower than what has been described in previous studies<sup>39,40</sup>. A lower number is expected, since unexpressed variants are not assayed by RNA-seq.

### Assaying isoform-specific genomic imprinting

To demonstrate the utility of the VG MPMAP–RPVG pipeline on a biological problem, we performed an exploratory analysis of genomic imprinting: a phenomenon in which some genes are expressed only from the copy inherited from a specific parent, regardless of its genomic sequence<sup>41</sup>. Several previous studies have studied imprinting genome-wide by quantifying ASE in RNA-seq data. These studies have demonstrated that imprinting varies across tissues<sup>41</sup> and varies in intensity across genes, with many genes showing biased expression but not complete silencing<sup>17,42</sup>. In addition, a handful of genes have been identified in which the polarity of imprinting depends on the isoform: some isoforms of the same gene are biased toward the paternal copy and others toward the maternal copy<sup>17</sup>.

The previous genome-wide studies have methodological limitations that diminish their ability to detect isoform-level imprinting. Some have aggregated ASE across all isoforms of the gene, which precludes isoform-level analysis<sup>41,42</sup>. The largest study, by Zink et al.<sup>17</sup>, performed tests on individual SNVs. This method can detect isoform-level differences in unshared exons. However, in shared exons, the ASE signal from the highest-expressed isoforms can drown out the signal of lower-expressed isoforms. Depending on the configuration of exons, this can make it very challenging to identify imprinting of opposite polarity.

Figure 5 shows results from our exploratory demonstration of isoform-level imprinting analysis using VG MPMAP and RPVG. We ran the entire pipeline using RNA-seq data from a lymphoblastoid cell line derived from 1000GP sample NA12878, which was sequenced as part of



**Fig. 5 | Exploratory demonstration of analyzing genomic imprinting using data from NA12878 lymphoblastoid cell line.** Results of the VG MPMAP–RPVG pipeline on RNA-seq data from a lymphoblastoid cell line from the ENCODE Project, focusing on genes previously identified as imprinted in blood. **a**, The proportion of expression attributed to the higher-expressed allele of heterozygous variants among the 20 most significantly imprinted genes from

the Zink et al. study<sup>17</sup> compared to all genes. The axes are scaled so that both histograms have the same area. **b**, Isoform-level haplotype-specific expression in *NAA60*, which was identified as imprinted but not as having isoform-dependent reversals in the polarity of imprinting in genome-wide studies. Isoforms with expression less than 0.25 TPM are not shown. Intervals indicate equal-tailed 90% credible intervals computed from 1,000 Gibbs samples.

the ENCODE project<sup>28</sup>. As a confirmatory analysis, we looked at the 20 ASE genes with the most significant *P* values from the Zink et al. study<sup>17</sup>. Mirroring that study's methods, we derived variant-specific ASE by summing over HSTs that contain a given allele. Figure 5a shows that the VG MPMAP–RPVG pipeline detects ASE at heterozygous variants in these imprinted genes at a markedly higher rate than in background across all genes.

The VG MPMAP–RPVG pipeline is also capable of detecting isoform-dependent genomic imprinting. Figure 5b shows an illustrative example in the gene *NAA60*. The isoforms show a complex pattern of imprinting polarity. Given the large differences in expression of these isoforms, the SNV-based analysis would have had difficulty identifying imprinting in the lowly expressed isoforms, and indeed this gene was reported as imprinted but not as having isoform-dependent imprinting<sup>17</sup>. However, this gene has been identified as having isoform-dependent imprinting using reverse transcription PCR in patients with uniparental disomies<sup>43</sup>. Nevertheless, it should be emphasized that this exploratory analysis, while suggestive, is insufficient to conclusively demonstrate isoform-dependent imprinting. Doing so would require further biological replicates and more rigorous controls for cis-regulation and cell line clonality<sup>42</sup>.

## Discussion

The pace of development in the field of eukaryotic pangenomics has surged in recent years. Improvements in sequencing technology have made it practical to characterize the genomes of increasingly many samples. As a result, pangenomes made from tens to hundreds of reference-quality genome assemblies have been constructed for many agricultural organisms<sup>44,45</sup>, and recently also for humans by the Human Pangenome Reference Consortium<sup>46</sup> and others<sup>47</sup>. Simultaneously, the bioinformatics tools to do pangenomic analyses have matured to the point of practicality for many applications<sup>9,48,49</sup>. We anticipate that pangenomic methods will continue to expand to inform increasingly many areas of genomics.

In this work, we have presented one step in this expansion: generalizing transcriptomics into pantranscriptomics. Our bioinformatics pipeline provides a full stack of tools for pantranscriptomic analysis. It can construct pantranscriptomes, map RNA-seq reads to these pantranscriptomes, and quantify transcription with haplotype resolution. The construction takes advantage of efficient pangenome data structures,

the mapping achieves a desirable balance of accuracy and speed, and the quantification can infer haplotype-specific transcript expression even when the haplotypes of the sample are not known beforehand.

Some downstream applications are already apparent. For one, the pipeline can be used to study causes of haplotype-specific differential expression. We demonstrated one such example by investigating genomic imprinting, uncovering suggestive evidence of complex patterns of imprinting at the isoform level. The pipeline could be similarly used to study other sources of haplotype-specific differential expression, such as nonsense-mediated decay and *cis*-regulation.

Another application is characterizing genotypes and haplotypes in coding regions from RNA-seq data. We demonstrated this capability by calling genotypes and HLA diplotypes. However, work is still needed to improve computational efficiency and accuracy in the HLA region. One of the major complications is that the dense variation in this region produces complicated graph topologies that lead to uncertainty in alignments.

For all of these applications, the VG MPMAP–RPVG pipeline increases the information that is available from RNA-seq data without paired genomic sequencing. This will enable low-cost study designs and deeper reanalyses of existing data.

The pipeline also has limitations. We have developed it to have good performance on pantranscriptomes constructed from phased variant calls. This is presently the most available data resource for constructing pangenomes. However, as increasingly many haplotype-resolved assemblies are produced, we predict that the emphasis in pangenomics will shift to pangenome graphs constructed from whole genome alignments. Constructing these graphs is currently an area of active research<sup>50,51</sup>. Such graphs have more complicated topologies. Experience suggests that pantranscriptomic tools will require further methods development to use these data resources effectively. This includes handling multi-mapping reads in RPVG, which will be crucial for inferring HST expression for genes in the more complex repetitive regions of these graphs.

Additional work on downstream analyses will be necessary to fully utilize HST expression inference. For example, current differential expression methods rely on comparing transcript counts between the same transcript of different individuals<sup>52</sup>. This is difficult at the HST level, since different individuals may not share a haplotype. While HST expression estimates can be marginalized to produce allele or



transcript expression estimates, more general statistical frameworks will need to be developed to avoid information loss between these steps in transcriptomic pipelines. A similar point holds for ASE estimation as well. Typical ASE pipelines include downstream statistical methods that assume known sample haplotypes. These do not readily accommodate haplotyping uncertainty that is inherent to the HST expression inference problem in RPVG.

Our pipeline is optimized for short read RNA-seq data. Long read RNA-seq technologies require specifically tailored algorithms for efficient analysis<sup>25</sup>. Pantranscriptomic analyses of long read RNA-seq data will likewise require further development. Nevertheless, the cost-effectiveness of short read sequencing ensures that it will remain an important part of the sequencing landscape into the near future.

Finally, our pipeline also relies on having a comprehensive pantranscriptome that contains many of the haplotype-specific transcripts from the sample. The pantranscriptomes used in this study (based on the 1000GP) provided good results in the three samples analyzed, but this performance may not extend to all other samples. Here—and throughout pangenomics—there is a compelling case to improve the completeness of data resources through more diverse sampling.

## Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41592-022-01731-9>.

## References

- Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-seq data with or without a reference genome. *BMC Bioinform.* **12**, 1–16 (2011).
- Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).
- Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods* **14**, 417–419 (2017).
- Degner, J. F. et al. Effect of read-mapping biases on detecting allele-specific expression from RNA-sequencing data. *Bioinformatics* **25**, 3207–3212 (2009).
- Eizenga, J. M. et al. Pangenome graphs. *Annu. Rev. Genomics Hum. Gen.* **21**, 139–162 (2020).
- Garrison, E. et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
- Rakocevic, G. et al. Fast and accurate genomic analyses using genome graphs. *Nat. Genetics* **51**, 354–362 (2019).
- Hickey, G. et al. Genotyping structural variants in pangenome graphs using the vg toolkit. *Genome Biol.* **21**, 1–17 (2020).
- Sibbesen, J. A., Maretty, L. & Krogh, A. Accurate genotyping across variant classes and lengths using variant graphs. *Nat. Genet.* **50**, 1054–1059 (2018).
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357 (2012).
- Rautiainen, M. et al. AERON: Transcript quantification and gene-fusion detection using long reads. Preprint at *bioRxiv* <https://doi.org/10.1101/2020.01.27.921338> (2020).
- Rautiainen, M. & Marschall, T. GraphAligner: rapid and versatile sequence-to-graph alignment. *Genome Biol.* **21**, 1–28 (2020).
- Denti, L. et al. ASGAL: aligning RNA-seq data to a splicing graph to detect novel alternative splicing events. *BMC Bioinform.* **19**, 1–21 (2018).
- Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
- Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory requirements. *Nat. Methods* **12**, 357–360 (2015).
- Zink, F. et al. Insights into imprinting from parent-of-origin phased methylomes and transcriptomes. *Nat. Genet.* **50**, 1542–1552 (2018).
- Castek, S. E., Levy-Moonshine, A., Mohammadi, P., Banks, E. & Lappalainen, T. Tools and best practices for data processing in allelic expression analysis. *Genome Biol.* **16**, 195 (2015).
- Van De Geijn, B., McVicker, G., Gilad, Y. & Pritchard, J. K. WASP: allele-specific software for robust molecular quantitative trait locus discovery. *Nat. Methods* **12**, 1061–1063 (2015).
- Rozowsky, J. et al. AlleleSeq: analysis of allele-specific expression and binding in a network framework. *Mol. Sys. Biol.* **7**, 522 (2011).
- Raghupathy, N. et al. Hierarchical analysis of RNA-seq reads improves the accuracy of allele-specific expression. *Bioinformatics* **34**, 2177–2184 (2018).
- Lee, W., Plant, K., Humburg, P. & Knight, J. C. AltHapAlignR: improved accuracy of RNA-seq analyses through the use of alternative haplotypes. *Bioinformatics* **34**, 2401–2408 (2018).
- Aguiar, V. R. C., César, J., Delaneau, O., Dermitzakis, E. T. & Meyer, D. Expression estimation and eQTL mapping for HLA genes with a personalized pipeline. *PLoS Genet.* **15**, e1008091 (2019).
- Sirén, J., Garrison, E., Novak, A. M., Paten, B. & Durbin, R. Haplotype-aware graph indexes. *Bioinformatics* **36**, 400–407 (2020).
- Wyman, D. et al. A technology-agnostic long-read analysis pipeline for transcriptome discovery and quantification. Preprint at *bioRxiv* <https://doi.org/10.1101/672931> (2020).
- Frankish, A. et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).
- Consortium, G. P. et al. A global reference for human genetic variation. *Nature* **526**, 68 (2015).
- Consortium, T. E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- Davis, C. A. et al. The encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res.* **46**, D794–D801 (2017).
- Berger, K., Somineni, H., Prince, J., Kugathasan, S. & Gibson, G. Altered splicing associated with the pathology of inflammatory bowel disease. *Hum. Genomics* **15**, 1–10 (2021).
- Micheletti, S. J. et al. Genetic consequences of the transatlantic slave trade in the Americas. *Am. J. Hum. Genet.* **107**, 265–277 (2020).
- Robinson, J. et al. IPD-IMGT/HLA database. *Nucleic Acids Res.* **48**, D948–D955 (2020).
- Chaisson, M. J. P. et al. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* **10**, 1784 (2019).
- Lappalainen, T. et al. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**, 506–511 (2013).
- Gourraud, P.-A. et al. HLA diversity in the 1000 Genomes dataset. *PLoS ONE* **9**, e97282 (2014).
- Abi-Rached, L. et al. Immune diversity sheds light on missing variation in worldwide genetic diversity panels. *PLoS ONE* **13**, e0206512 (2018).
- Orenbuch, R. et al. arcasHLA: high-resolution HLA typing from RNAseq. *Bioinformatics* **36**, 33–40 (2019).
- McLaren, W. et al. The Ensembl variant effect predictor. *Genome Biol.* **17**, 122 (2016).
- Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).



40. Maretty, L. et al. Sequencing and de novo assembly of 150 genomes from Denmark as a population reference. *Nature* **548**, 87–91 (2017).
  41. Baran, Y. et al. The landscape of genomic imprinting across diverse adult human tissues. *Genome Res.* **25**, 927–936 (2015).
  42. Jadhav, B. et al. RNA-seq in 296 phased trios provides a high-resolution map of genomic imprinting. *BMC Biol.* **17**, 1–20 (2019).
  43. Nakabayashi, K. et al. Methylation screening of reciprocal genome-wide UPDs identifies novel human-specific imprinted genes. *Hum. Mol. Genet.* **20**, 3188–3197 (2011).
  44. Liu, Y. et al. Pan-genome of wild and cultivated soybeans. *Cell* **182**, 162–176 (2020).
  45. Zhou, Y. et al. Graph pangenome captures missing heritability and empowers tomato breeding. *Nature* **606**, 527–534 (2022).
  46. Liao, W.-W. et al. A draft human pangenome reference. Preprint at *bioRxiv* <https://doi.org/10.1101/2022.07.09.499321> (2022).
  47. Ebert, P. et al. Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021).
  48. Sirén, J. et al. Pangenomics enables genotyping of known structural variants in 5202 diverse genomes. *Science* **374**, abg8871 (2021).
  49. Ebler, J. et al. Pangenome-based genome inference allows efficient and accurate genotyping across a wide spectrum of variant classes. *Nat. Genet.* **54**, 518–525 (2022).
  50. Li, H., Feng, X. & Chu, C. The design and construction of reference pangenome graphs with minigraph. *Genome Biol.* **21**, 1–19 (2020).
  51. Hickey, G. et al. Pangenome graph construction from genome alignment with Minigraph-Cactus. Preprint at *bioRxiv* <https://doi.org/10.1101/2022.10.06.511217> (2022).
  52. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 1–21 (2014).
- Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.
- Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.
- © The Author(s), under exclusive licence to Springer Nature America, Inc. 2023

## Methods

### Sequencing data, transcript annotations, and variation data-bases

GENCODE v29 (primary assembly) was used as a transcript annotation set<sup>26</sup>. All transcripts with either the mRNA\_start\_NF or mRNA\_end\_NF tag were removed to only keep confirmed full-length transcripts. Furthermore, a transcript subset containing 80% of the GENCODE transcripts was created by randomly removing 34,490 of the 172,449 transcripts in the annotation. The fraction removed was based on recent estimates of the fraction of novel transcripts in a sample using long reads<sup>25</sup>.

Genomic variants on GRCh38 from the 1000GP were downloaded from EBI ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/GRCh38\\_positions/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/GRCh38_positions/))<sup>27</sup>. The variants were first normalized using BCFTOOLS<sup>53</sup> and four different sets containing variants from differently-sized collections of samples were created (Supplementary Table 1). Two of these sets were constructed so as to not include variants unique to the CEU population. This was because we benchmarked the pipeline on NA12878, who is from this population, and we wanted our evaluations to approximate an expected use case of sequencing a sample from a population that is not represented in the reference haplotype panel. For all of the variant sets, the intronic and intergenic variants were further filtered using BCFTOOLS, keeping only variants with an alternative allele frequency of at least 0.002 or 0.001 depending on the set. This was done to decrease the complexity of the graph in regions where fewer reads are expected to map. The GRCh38 (primary assembly) reference genome used throughout the study was downloaded from Ensembl ([http://ftp.ensembl.org/pub/release-94/fasta/homo\\_sapiens/dna/](http://ftp.ensembl.org/pub/release-94/fasta/homo_sapiens/dna/)).

A list of all sequencing data used can be found in Supplementary Table 4.

### Spliced pangenome graph construction

We developed a method in the VG toolkit, VG RNA, for constructing spliced pangenome graphs from a transcript annotation and an existing pangenome graph. VG RNA begins by identifying the path in the graph that corresponds to each exon in the annotation. These exon paths can start or end internally in a node rather than only at boundaries between nodes, as with other paths in VG. Next, VG RNA divides nodes as necessary to expose exon boundaries as node boundaries and then adds edges (splice junctions) to the graph connecting adjacent exons within each transcript. The transcript paths are then labeled in the resulting spliced pangenome graph. Lastly, the node ID space of the spliced pangenome graph is compacted and reordered in topological order to make graph compression more efficient<sup>54</sup>.

Different combinations of transcript annotations (full and an 80% random subset) and variant sets (Supplementary Table 1) were used as input to create the graphs used in the mapping and expression inference evaluation (Supplementary Table 2).

### Pantranscriptome construction

Alongside spliced pangenome graphs, VG RNA can simultaneously generate pantranscriptomes consisting of HSTs created from transcript and haplotype annotations. It creates pantranscriptomes by projecting the reference transcript paths onto haplotypes paths indexed using the GBWT<sup>24</sup>, a data structure for efficiently storing thousands of paths in a graph, such as haplotypes or transcripts. If nodes are split during the spliced pangenome graph construction (see above), VG RNA first updates the haplotypes in the input GBWT. Next, the flanking positions of the exon boundaries on the reference chromosome path are used as anchors for projecting exons between the reference and haplotype paths. Anchoring on the positions adjacent to exon boundaries allows for genomic variation at the distal ends of exons.

To find all possible haplotype paths between two exon anchors, we use an exhaustive depth-first search (DFS) initialized at the start anchor.

Branches in the DFS (branch) are queried against the GBWT index and terminated if they are not a subpath of any haplotype. A search is also terminated if none of the haplotypes in the branch contain the end anchor node. The output from the search is a list of unique projected haplotype-specific (HS) exon paths and the set of haplotypes consistent with each of them. The final HST paths are constructed one exon at a time by connecting HS exon paths that share at least one haplotype for each transcript. The fact that all the HS exon paths are unique makes the approach scale well with the number of haplotypes, as it can take advantage of the fact that haplotypes are often identical locally.

A list of all pantranscriptomes created for this study including the transcript annotations and variant sets used as input can be seen in Supplementary Table 3. The HSTs were written both as nucleotide sequences in FASTA format (for inputs to other expression inference tools) and as paths to a GBWT. A bidirectional GBWT, where each path is stored in both directions, was also created. RPVG uses this index to decrease computation time when reads are not strand specific. For each GBWT, a corresponding *r*-index was constructed, which decreases the computation time it takes to query path IDs in the GBWT<sup>55</sup>.

### Read simulation model

Most of the simulated reads were generated using VG SIM, a read simulator in the VG toolkit that is designed primarily for next-generation sequencing reads. Its model consists of three components: a Markov model for base quality strings, a path frequency model, and a fragment length model (when sampling paired-end reads).

The model for base quality strings is fit to replicate the base quality strings in a user-provided FASTQ. A separate Markov transition distribution is fit for each base position in the read. The state of each Markov distribution consists of two components: the Phred base quality at that base and whether that base is an N. If a paired-end FASTQ is provided, VG SIM will fit a separate model for each read end. In addition, the first states of each read in the pair are modeled with a single joint distribution, which allows for some dependence between the quality of both reads in the pair. The probabilities of the Markov transitions and the initial states are estimated by their empirical frequency in the FASTQ.

VG SIM determines the base sequence of each read by following random walks through the pangenome graph. These walks may optionally be restricted to specific paths through the graph, such as paths of transcripts in a spliced pangenome graph. The sampling frequency of a transcript path is proportional to the product of its effective length<sup>4</sup> and its expression value measured in TPM, as determined by a user-provided expression profile. Once the path has been chosen, the starting location of the read is selected uniformly at random along the transcript. The sequence of the walk is then extracted, and sequencing errors are introduced according to the probability distribution implied by the base quality string. A user-specified fraction of these errors are produced as indel errors rather than substitution errors.

When simulating paired-end sequencing, the fragment length is modeled with a normal distribution. The user provides the mean and standard deviation for this distribution. For a given path, the normal distribution is truncated to between one and the path length. Both reads are sampled from a single walk through the graph with length equal to the sampled fragment length. If this length is shorter than the read length, the read is truncated to the fragment length.

### Simulating RNA-seq reads from haplotype-specific transcripts

Reads were simulated from HST paths derived from the haplotypes of NA12878 in the 1000 Genomes Project (1000GP) and the GENCODE transcript annotation. The corresponding spliced pangenome graph was created using VG RNA.

In total, we created five different simulated read sets: four using VG SIM and one using RSEM<sup>1</sup>. Two different read sets were used to fit the simulations' error model: SRR1153470 and ENCSTR000AED, replicate 1. For both real read sets, we used VG SIM to create two simulated read

sets. One set of reads was simulated with an expression profile derived from the real data, and the other set was simulated with uniform expression across transcripts. The single RSEM simulation used the uniform approach. Supplementary Table 5 lists all of the simulated read sets. Each was used in different parts of the benchmarking. The uniform expression datasets were used to benchmark mapping whereas the data-based expression read sets were used to benchmark expression inference. We used a uniform expression profile for the mapping benchmark to not bias the analyses towards easily-mappable transcripts.

To ensure balanced expression between the two haplotypes for all transcripts, only transcripts that were successfully projected to both haplotypes were given a positive expression for the uniform expression set. For the simulated read sets with data-derived expression, we generated the expression profile by mapping the reads using BOWTIE2<sup>11</sup> with default parameters and then quantifying using RSEM, also with default parameters. For all five read sets, we simulated 25 million 101 base-pair read pairs from each haplotype. For VG SIM, we used an indel probability error of 0.001 and the base quality distribution was trained on 10 million randomly sampled read pairs of the training data. The read pairs were sampled using seqtk (<https://github.com/lh3/seqtk>). RSEM was given the estimated training data model file and the background noise fraction was set to zero.

### Mapping and multipath alignment with VG MPMAP

Like most read mappers, the mapping algorithm of VG MPMAP is designed using the ‘seed–cluster–extend’ paradigm. First, it locates exact matches, ‘seeds’, between the read and the graph. Next, the seeds are ‘clustered’ together to identify regions of the graph that the read could align to. Finally, the seeds are ‘extended’ into an alignment of the entire read. Because these operations occur in the context of a pangenome graph, they use several specialized algorithms and indexes.

**Seeding.** VG MPMAP seeds alignments with maximal exact matches (MEMs) against the graph, which it finds using a GCSA2 index<sup>56</sup>. MEMs are exact matches between an interval of the read and a walk in the graph such that the match cannot be extended further in either direction at that location in the graph. The MEMs are found using a two-stage algorithm, which has also been described previously<sup>7</sup>.

In the first stage, the algorithm finds super-maximal exact matches (SMEMs), which are MEMs for which the read interval is not contained within the read interval of any other MEM (Supplementary Algorithm 1). This algorithm also relies on a longest common prefix array. The second stage of the algorithm finds the longest MEMs that are shorter than each SMEM but have their read interval contained in the read interval of the SMEM, subject to a minimum length (Supplementary Algorithm 3).

**Clustering.** The clustering algorithm in VG MPMAP is built around the distance index described in Chang, et al.<sup>57</sup>. In brief, this index can query the minimum distance between two positions in the pangenome graph by expressing the distance as the sum of a small number of precomputed distances. This is accomplished by taking advantage of the common topological features of pangenome graphs, namely that they tend to contain long chains of bubble-like motifs that result from genomic variation<sup>58</sup>.

The clustering algorithm begins by constructing a directed acyclic graph (DAG) in which the nodes correspond to MEM seeds. The edges are added whenever (i) there is a path connecting two seeds in the graph, and (ii) the seeds are collinear along the read. We use the distance index to determine the existence of a path that connects the seeds in the graph, and the edges are also labeled by the distance. Edges that are much longer than the read length are not added; this avoids treating distal elements on the same chromosome as part of the same cluster. In addition, we accelerate this process using Algorithm 3 from Chang et al.<sup>57</sup>, which partitions seeds into equivalence classes on the basis of the distance between them. The equivalence relation is the

transitive closure of the relation of being connected by a path of length less than  $d$ , which is a tunable parameter. By choosing  $d$  correctly, we can ensure that all of the edges we would include occur between seeds in the same equivalence class. This significantly reduces the number of distance queries we need to perform.

Once the DAG of seeds has been constructed, we approximate the contribution of each seed and edge to the score of an alignment that contains them. Seeds are scored as a short alignment of matches, and edges between seeds may be scored as an insertion or deletion if the distance in the graph does not match the distance on the read. We then use dynamic programming to compute the heaviest path defined by the node and edge weights (scores) within each connected component and take the seeds along this path as a candidate cluster. Clusters are passed through to the next stage of the algorithm if their weight is within a prespecified amount of the heaviest-weight cluster, subject to a hard limit on the total number of clusters.

**Multipath alignments.** Most existing sequence-to-graph aligners, including VG MAP<sup>7</sup>, produce an alignment of the sequence to a particular walk through the graph. VG MPMAP uses a different alignment formalism, which we call a multipath alignment. In a multipath alignment, the sequence can diverge and reconverge along different walks through the graph (Extended Data Figure 1). Thus, the read can align to a full subgraph rather than to a single path. This allows the alignment object to carry within itself the alignment uncertainty at known variants or splice junctions. This information can be used in downstream inference applications, including RPVG.

More formally, a multipath alignment of read  $R$  is itself a digraph with the following properties:

1. Each node corresponds to an alignment of some substring of  $R$  to a path in the pangenome
2. An edge between  $u$  and  $v$  exists only if  $u$  and  $v$  align adjacent substrings of  $R$  to adjacent paths in the pangenome.
3. Every source-to-sink path through the multipath alignment can be concatenated into a complete, valid alignment of  $R$  to a path in the pangenome.

VG MPMAP additionally annotates the partial alignment of each node with its alignment score. The alignment score of any particular sequence-to-path alignment expressed in the multipath alignment can be computed efficiently by simply adding the partial alignment scores along the path.

While sequence alignments have well-established optimization criteria, there is no such criterion for optimizing the topology of a multipath alignment. In lieu of one, we adopt heuristics that are motivated by the common topological features of pangenome graphs. Our high-level strategy is to use exact match seeds to anchor alignments. We then align between seeds and within sites of variation in the graph.

**Anchoring alignments.** To use a cluster of exact match seeds to anchor a multipath alignment, it is first necessary to compute the reachability relationships between the seeds. This is a non-trivial problem.

We begin by converting the local graph around a cluster into a directed acyclic graph using an algorithm that has been described previously<sup>7</sup>. In brief, we identify small feedback arc sets within each strongly connected component using the Eades–Lin–Smyth algorithm<sup>59</sup>, and then we duplicate the strongly connected components with the feedback arcs linking successive copies. Using dynamic programming over the DAG as we construct it, we can preserve all cyclic walks up to some prespecified length, which is based on the read length.

After creating the DAG, we inject the seeds into the new graph. Since the DAG conversion algorithm can expand the node space of the original graph, seeds can now correspond to multiple locations in the DAG. In this case, we duplicate the seeds to all of the corresponding

locations in the DAG. We then use a three-stage algorithm that computes the transitive reduction of a graph in which the nodes correspond to seeds, and two seeds have an edge between them if they are collinear along the read and reachable within the pangenome graph (Supplementary Algorithm 4).

1. Compute the reachability relationships between the seeds, ignoring collinearity on the read.
2. Rewire the reachability edges between the seeds to respect collinearity on the read.
3. Compute the transitive reduction of the resulting graph.

This algorithm is designed to have linear run time in the number of seeds and the size of the DAG, but only in the typical case where the seeds line up along a walk through the pangenome graph. In the general case, the run time can be quadratic.

**Dynamic programming with multiple traceback.** The alignments between anchors are computed using a banded implementation of partial order alignment<sup>60</sup>. The alignments of the read tails past the end of anchors are computed using a SIMD-accelerated POA implementation from the gssw library (<https://github.com/vgteam/gssw>).

We use a specialized traceback algorithm to obtain the alignments to multiple paths through the pangenome graph from a single dynamic programming problem (Supplementary Algorithm 8). The algorithm returns the  $k$ -highest-scoring alignments, and we choose  $k$  to be the number of paths through the subgraph we are aligning to, subject to a hard maximum. The key insight behind the algorithm is that the next highest-scoring traceback can be determined by checking local properties of the dynamic programming matrix while computing the highest-scoring traceback. In addition, for each anchor that crosses a snarl (bubble-like graph features that often indicate variation<sup>38</sup>), we remove the interior of snarl before performing alignments. This way, the multiple traceback algorithm can align to multiple paths at sites of variation.

**Quantifying mapping uncertainty.** The method that VG MPMAP uses to compute mapping quality is largely shared with VG MAP (see Garrison, et al.<sup>7</sup>; Supplementary Note). As in VG MAP, base qualities are incorporated into alignment scores (essentially downweighting low-quality bases), and the alignment scores are subsequently used to compute a mapping quality. The formulas used to compute mapping quality rely on the conversion of alignment scores into the log-likelihood of a hidden Markov model.

VG MPMAP also uses a concept of the ‘multiplicity’ of a mapping to model errors introduced by the mapping algorithm itself. In particular, at certain points in the algorithm, we enforce hard caps on certain algorithmic behaviors, such as the number of alignments that will be attempted. If we run up against these hard caps, we expect that not all high-scoring alignments will be found. We incorporate this information into the mapping quality formula by treating alignments as if multiple equivalent alignments actually were found. For example, if we attempted alignments for 10 of 30 promising clusters and found 1 high-scoring alignment, we would estimate its multiplicity to be 3. We then compute the mapping quality as if two additional copies of the alignment had been found.

Multiplicities allow VG MPMAP to aggregate information about sources of algorithmic inaccuracy over different steps in the algorithm. The central entities in each step of the mapping algorithm (seeds, clusters, alignments, and read pairs) are each associated with a multiplicity. When combining orthogonal pieces of information (seeds in a cluster, or single-end alignments in a paired alignment), the new entity receives the minimum of its constituents’ multiplicities. When layering on a new source of algorithmic uncertainty (a further hard cap), the multiplicity of an entity is multiplied by its estimated multiplicity in that step of the algorithm.

**Determining statistical significance.** VG MPMAP uses a frequentist hypothesis test to assess the statistical significance of a read alignment. The test statistic that we use is the alignment score. The null hypothesis is that the alignment score was obtained by a uniform random sequence of the same length as the read. By default, we set the type-I error rate to 0.0001. If the  $P$  value of an alignment score is not significant at this level, the read is reported as unmapped.

Modeling the null hypothesis of the test is not entirely straightforward. In general, we expect higher local alignment scores from longer reads or larger pangenome graphs. However, there are subtleties. A large pangenome graph may consist of many repeats of the same sequence so that its effective size is smaller than its total sequence length. Alternatively, a small graph may have a complex topology that admits a combinatorially large set of walks. For these reasons, we take an empirical approach that fits a model to match the pangenome graph. At the start of every mapping run, we map a sample of uniform random sequences of varying lengths and use the scores to fit the parameters of a distribution using maximum likelihood. Those parameters are then regressed against the read length. The regression allows us to query the  $P$  value for a read of any length.

The parametric distribution we use can be derived as the maximum of  $v$  independent, identically distributed exponential variables with rate  $\lambda$ , which has the following probability density function:

$$f(x|\lambda, v) = \lambda v (1 - e^{-\lambda x})^{v-1} e^{-\lambda x}. \quad (1)$$

The fitting algorithm alternates between maximizing the likelihood with respect to each of the two parameters with the other fixed until convergence.  $v$  is fit using the Newton-Raphson method, and  $\lambda$  is fit using golden-section search.

The motivation for this model is that the length of the match starting at position of a uniform random sequence (the read) and position of a fixed sequence (the reference) is approximately Geometric(1/4). The optimal local alignment score is closely related to the longest match at any position on the read sequence to any position on the pangenome graph. We use an exponential distribution because it closely approximates a geometric distribution and is easier to fit.

**Paired-end mapping.** At the beginning of each paired-end mapping run, VG MPMAP uses a sample of the first 3,000 uniquely mapped pairs to fit parameters of a fragment length distribution. The distance between the reads in each pair is computed with the distance index. Non-uniquely mapped pairs are buffered and then remapped after the fragment length distribution has been fit.

The fragment distribution is modeled as a normal random variable with mean  $\mu$  and variance  $\sigma^2$ . We use a method of moments estimator for a truncated normal distribution so that the parameter estimation is robust to erroneous distances from possible mismappings or unannotated splice junctions. In particular, we discard the largest and smallest  $\frac{1-\gamma}{2}$  fraction of fragment length measurements (default  $\gamma = 0.95$ ). The remaining  $\gamma$  fraction of measurements correspond to a sample from a truncated normal distribution with the same  $\mu$  and  $\sigma^2$ . The following estimators can be derived using method of moments on this truncated normal distribution:

$$\begin{aligned} \hat{\mu} &= \bar{x} \\ \hat{\sigma}^2 &= s^2 \left( 1 - \frac{2\phi(\alpha)}{\gamma} \right)^{-1}, \end{aligned} \quad (2)$$

where  $\phi$  is the density function of a standard normal distribution,  $\bar{x}$  and  $s^2$  are the empirical mean and variance among the retained measurements, and  $\alpha = \Phi^{-1}(\frac{1-\gamma}{2})$  is the left truncation point on a standard normal distribution.

When mapping paired-end reads, the clustering stage of the algorithm adds an additional step. First, each read in the seeds of pair are



clustered as in the single-end algorithm. Next, the clusters from the two reads are paired by checking which pairs imply a fragment length within ten standard deviations of the mean, as estimated by the algorithm in the previous section. The implied fragment length connecting two clusters is estimated using the distance index, with the position of a cluster taken to be the position of its longest seeds. Pairs of clusters are prioritized by a sum of an estimated alignment score (interpreted as a log-likelihood) and the log-likelihood of the normal distribution that we model the fragment length distribution with.

It sometimes happens that the mapping heuristics fail on only one of the two reads of a fragment. When this occurs, it is sometimes possible to ‘rescue’ the alignment of the other read by aligning it to the region of the pangenome graph where we expect to find it relative to the mapped read. VG MPMAP employs this strategy whenever the pair clustering procedure fails to produce a pair of clusters consistent with the fragment length distribution, or when all of the clustered alignment pairs have at least one end without a statistically significant alignment. We also perform a limited number of rescues even when a consistent cluster pair is found, if there are clusters of at least one of the ends that are equally as promising as the one in the cluster pair. We place a hard cap on the number of rescues performed to control run time. The fraction of eligible rescues that were actually performed becomes a component in the multiplicity of an alignment.

The multipath alignment algorithm is slightly different when computing rescue alignments, because there are no exact match seeds to use as anchors. Instead, we first perform a single-path alignment using gssw. Then we remove any sections of the alignment that lie inside snarls, and realign those segments of the read as when connecting anchors in the standard multipath alignment algorithm.

**Spliced alignment.** Because spliced pangenome graphs include annotated splicing events as edges, it is usually unnecessary to use specialized alignment algorithms to obtain spliced alignments. However, transcript annotations are incomplete, so it is still important to be able to produce spliced alignments. VG MPMAP includes a spliced alignment algorithm but applies it conservatively: only when an alignment includes a moderately long soft-clip on at least one end. A long soft-clip is suggestive that the clipped end of the read might align to a part of the graph that was too distant to be included in the primary seed cluster, as would be expected with an unannotated splice event. One common exception to this pattern are adapter sequences, which can be captured in a read when the sequenced fragment is shorter than the read length. To avoid the computational burden of attempting to find nonexistent spliced alignments for these cases, common adapter sequences are specifically excluded from this subroutine.

The spliced alignment algorithm begins by finding candidate regions to align the clipped read end to. These regions are selected by scanning over secondary mappings, unaligned seed clusters, and unclustered seeds. For paired reads, spliced alignments can also be found by rescuing the soft-clipped portion of the read from the other read in the pair. This is only possible when the soft-clip is on the side of the read that faces inward on the fragment. Spliced alignment rescue is only attempted when none of the other spliced alignment candidates yields a statistically significant spliced alignment.

Spliced alignment candidates must pass several filters to be included in the read mapping. Candidates must roughly correspond to the clipped end of the read. They must also be reachable from the primary alignment by some path in the graph, which is determined using the distance index. The final filter that a spliced alignment must pass is a significance test. This test has three components: (i) the increase in alignment score that results from aligning the additional bases; (ii) the bias against the splice site motifs in the intron; and (iii) the bias against the intron length. The default parameterization is based on the human transcriptome. Splice site motifs are penalized by their log-frequency, as given by Burset et al.<sup>61</sup>. The bias against the intron

length is determined using the log-likelihood of a log-normal mixture model fit to the human intron length distribution. The three components are combined into a joint log-likelihood and tested against a critical value.

To compute the quantities needed for this test, the spliced alignment algorithm identifies the splice motifs near the ends of a pair of splice candidates. If any pair of canonical splice site dinucleotides are found on any path from the two ends, the intervening sequence is aligned as if the two splice sites were joined by an edge in the graph. In addition, the intron length is measured between these positions in the graph using distance along the reference path.

**Multi-mapping reads.** Reads with multiple high-scoring alignments can be reported in two different ways. First, separate alignments can be reported up to a user-specified maximum number (default 10). Second, a single, possibly disconnected multipath alignment can be reported that includes all high-scoring alignments. In the first option, all of the reads are annotated with a collective ‘group mapping quality’ that quantifies the probability that all of the reported alignments are incorrect. In the latter option, the main mapping quality annotation is equivalent to the group mapping quality.

### RNA-seq mapping evaluation

We compared the performance of VG MPMAP at mapping RNA-seq data against the existing graph alignment method VG MAP<sup>7</sup> from the VG toolkit and two state-of-the-art RNA-seq mapping tools, HISAT2<sup>15</sup> and STAR<sup>2</sup>. Graph indexes and genomes were created for each tool using default parameters, with MPMAP and MAP sharing the XG and GCSA index. All mappers were run with default or recommended parameters for RNA-seq data. For the simulated data the maximum number of reported multi-alignments per read was set to ten for each method.

The main mapping results were obtained using the ENCSR000AED, replicate 1 data (see Supplementary Table 4): Fig. 2; Extended Data Figs 2–5, and Supplementary Figs 1–4 and 6. The SRR1153470 and CHM13 data (see Supplementary Table 4) were used to optimize the parameters of VG MAP and VG MPMAP. Nevertheless, the pattern of performance on these reads is similar to the ENCSR000AED data (Supplementary Figures 5 and 7).

We evaluated mapping accuracy on simulated reads using two different methodologies to ensure the robustness of our conclusions. One methodology was based on basewise overlaps along the linear reference genome, and the other was based on distances along transcript and reference paths in the graph. In both cases, the results were stratified by mapping quality. For VG MPMAP, we used the group mapping quality (see above), and for the other tools we used the mapping quality value of the alignment with the highest overlap or closest distance, or the highest of these mapping qualities in case of ties.

For the overlap-based evaluation, the graph alignments were first projected to the reference paths using VG SURJECT in spliced alignment mode. In brief, VG SURJECT takes a set of graph-aligned reads and realigns them to all nearby reference paths in the graph, producing a BAM file with the reads aligned to the reference sequences. The re-alignment is only performed on the parts of the alignment that do not already follow the reference paths. A read was considered correctly mapped if 90% of the bases of the simulated true reference alignment were covered by one of its multi-alignments. The true reference alignments were generated using the transcript position of each read provided by VG SIM or RSEM, and the NA12878 haplotype-specific transcript reference alignments. The latter were created by projecting the transcript paths to the reference sequences using VG SURJECT in spliced alignment mode.

Owing to sequencing artifacts, the ends of reads will occasionally consist of such low-quality bases as to be practically random. Our simulation framework recapitulates this feature of real sequencing data. However, in real data these read ends do not correspond to any

underlying genomic sequence, whereas the simulation assigns them a true genomic alignment. Aligners that soft-clip these uninformative bases would be penalized in this evaluation, even though this is the correct decision for real data. We therefore trimmed all bases at both ends of an alignment (including the true alignments) that had a Phred base quality score below 3 for the purpose of computing the overlap. All alignments for which more than half of the sequence was trimmed were discarded from the evaluation so that the percent overlap could be estimated more confidently.

To classify whether a read contained any novel (unannotated) splice junctions we looked at all deletions and reference skips in the true reference alignment with a length of at least 20 base pairs. These were compared to the transcript annotation that was used to build the graph or reference, and defined as novel if it was not possible to find a splice junction in the annotation that was within 5 base pairs at both ends.

Edit distance was calculated as the number of base pair differences between the simulated read and the reference sequence. The NA12878 1000GP genotypes were used to estimate the genomic variation edit distances.

We used the VG GAMPCOMPARE tool for the distance-based evaluation. The truth set in this evaluation was the true graph alignments produced by VG SIM. In short, VG GAMPCOMPARE finds the minimum possible distance between the start position of an estimated alignment and the true alignment across all reference and transcript paths in the graph. Before running VG GAMPCOMPARE, the BAM format alignments from HISAT2 and STAR were converted into graph alignments (GAM format) using VG INJECT, which translates linear reference alignments into alignments against the path of the reference in a graph. An alignment was considered correct if its start position was within 100 base pairs of the start position of the true alignment along the path of the reference or any transcript path.

Reference bias was quantified using simulated reads, by counting the number of reads with a mapping quality value of at least 30 that overlapped heterozygous variants. For this analysis, we used the linear reference-based alignments. To treat different variant types and lengths equally, we computed the read count for each variant allele as the average read count across the two breakpoints of the allele. Reads simulated from each haplotype were counted separately and only variants with at least 20 reads across both alleles combined were used to quantify reference bias. We skipped variants that were not classified as SNVs, simple deletions, or simple insertions.

To further evaluate allelic bias, we counted the reads supporting each allele of heterozygous variants among mapped simulated reads for each of the mappers using the approach described above. We also added a pipeline consisting of STAR followed by read filtering with WASP to this comparison. We found that WASP was computationally infeasible using the full CEU population, so we instead gave it a variant database consisting of only variants from NA12878. Therefore, to have a better comparison to WASP, we also created sample-specific references for VG MPMAP, VG MAP, and HISAT2, and we report results for these references as well. We estimated the observed rate of false positives by testing for allelic skew in mapped reads on heterozygous variants using a two-sided binomial test ( $\alpha = 0.01$ ). All significant *P* values are false positives, since the reads were simulated without an allelic bias.

When benchmarking using real reads, truth alignments are not available. Instead, we used a proxy measure of aggregate mapping accuracy on the basis of long read mappings from the same cell line. The long reads are easier to map confidently, and we expect the cell line to have similar transcript expression across replicates. Thus, higher correlation between the coverage of short read mappings and the coverage of long read mappings is indirect evidence of higher accuracy. For long read data, we used NA12878 PacBio Iso-Seq alignments generated by the ENCODE project (Supplementary Table 4). The cleaned Iso-Seq alignments of four replicates were first merged and secondary alignments and alignments with a quality below 30 were filtered using

SAMTOOLS<sup>62</sup>. These filtered alignments were then compared to the short read RNA-seq alignments by calculating the Pearson's correlation of the average exon read coverage between the two. Exons were defined using the Iso-Seq alignments by first converting them to BED format and then merging overlapping regions using BEDTOOLS<sup>63</sup>.

We measured memory and compute time for all mappers using the Unix time utility. The mapping compute usage of each tool was estimated using 16 threads on an m5.4xlarge Amazon Web Services (AWS) instance. The reads per second statistic was computed by dividing the number of reads by the product of the wall clock time and the number of threads. This is a somewhat biased measurement, since it includes the one-time start up computation that does not scale with the number of reads. However, the magnitude of this bias is small, and it tends to disfavor VG MPMAP, which has the longest start up of the tools we evaluated.

Reference alignments in BAM format were sorted and indexed, also using SAMTOOLS. The SeqLib library was used in the evaluation scripts to parse the alignments and calculate overlaps<sup>64</sup>.

### Haplotype-specific transcript quantification

We developed RPVG as a general tool for inferring the most likely paths and their abundance from a set of mapped sequencing reads. In this study we used RPVG to quantify the expression of haplotype-specific transcripts (HSTs) in a pantranscriptome. The RPVG algorithm consists of four main steps:

1. Find read alignment paths that align to HST paths.
2. Cluster alignment paths and HST paths.
3. Calculate alignment path probabilities.
4. Infer haplotypes and expression from probabilities.

A graphical overview can be seen in Supplementary Fig. 18.

**Finding alignment paths.** The first step of RPVG is to parse each alignment and find all alignment paths that follow at least one HST path in the pantranscriptome GBWT index (Supplementary Figure 18a). An alignment path is the set of nodes a read alignment follows in the graph. For single-path alignments there is only one alignment path, but for multipath alignments there can be many. We will focus here on multipath alignments, since a single-path alignment is a sub-case when a multipath alignment only contains a single path.

Multipath alignments are represented as a graph, and thus the objective is to find all paths through this graph that also exist as subpaths in the GBWT. This search would normally scale linearly in the number of HSTs overlapping the read, but the GBWT allows us to simultaneously query all locally identical HSTs that contain the same subpath.

RPVG uses a depth-first search (DFS) through the multipath alignment graph to find all alignment paths. A branch in the search is terminated if its alignment path is not present as a subpath in the GBWT. A DFS is initialized at each source node in the alignment graph. We terminate any alignment path early where it is not possible to reach a score of 20 below the current highest-scoring path, assuming perfect scoring for the remainder of the alignment.

The topology of the multipath alignment graphs is determined by heuristics. In some cases these heuristics fail, resulting in multipath alignments that do not cover all possible alignment paths. This can result in incorrect downstream expression estimates as a read might be missing an alignment path to the correct HST. To overcome this, RPVG allows alignment paths to be shortened in order to be made consistent with an HST path. More specifically, the DFS can start and end up to four bases inside the read (excluding soft-clipped bases). The score of partial alignment paths are penalized proportionally to the number of non-matched bases at each end, adjusted for their quality. The longest possible alignment path to an HST is selected as the best alignment.

The output from the DFS is one set of alignment paths for each multipath alignment. Next, RPVG labels a set as low scoring if the highest-scoring alignment path in the set is less than 0.9 times the maximum possible quality-adjusted alignment score. The sets labeled as low scoring are treated as being incorrect; they may be misalignments, or they may originate from an HST not in the input pantranscriptome. They are later used when calculating the noise probability.

For paired-end reads, one additional step is needed: combining the alignment paths of each read to create a set of alignment paths for the whole fragment. First a set of alignment paths is generated for each alignment in the pair as described above. Next, RPVG attempts to combine each start (first read) alignment path with each of the end (second read) alignment paths. If the fragments are not strand-specific and the pantranscriptome GBWT is not bidirectional, RPVG then repeats the process using the reverse complement of the fragment.

The procedure to combine the two alignment paths differs depending on whether they overlap or not. If they do overlap, a single combined alignment path is created for the fragment by merging the two while requiring that the path of overlapping portions matches perfectly. If they are separated by an insert, the start alignment path is extended using a depth-first search following the HST paths. If the search reaches one of the start nodes for an end alignment path, a new fragment alignment path is created by merging the search and end alignment path. The new fragment alignment path is only kept if it follows at least one HST path in the pantranscriptome. The search is terminated if all start nodes in the end alignment paths have been visited and they are not part of a cycle. An alignment path is discarded if its length is above  $\mu + 10\sigma$ , where  $\mu$  and  $\sigma$  are the mean and standard deviation of the fragment length distribution. These parameters are either supplied by the user or parsed from the input alignments (the VG aligners write the parameters they estimated to the alignment file). The score of the resulting fragment alignment path is calculated as the sum of the scores of the two read alignment paths. The mapping quality is calculated as the minimum across the two reads.

The final output from the search is a set of alignment paths and the HSTs that each path aligns to for each read or fragment. For simplicity, in the following, we will use the term “fragment” to denote both a single-end read and a set of paired-end reads.

**Clustering transcript paths.** HST paths that do not share any fragments are independent, and therefore their expression can be inferred separately. By contrast, the expression of HST paths that share alignments must be inferred jointly. Accordingly, RPVG identifies clusters of HST paths that share alignment paths from the same fragment. By dividing the inference problem into these smaller, independent clusters, computation and memory can be considerably reduced. The clustering algorithm works by first constructing an undirected graph where vertices correspond to HST paths and edges correspond to HST paths being observed in the same set of fragment alignment paths. Connected components in this graph correspond to clusters.

**Calculating alignment path probabilities.** For each fragment, the probability of it originating from each of the HSTs in its cluster is calculated by RPVG using the alignment path scores, lengths and mapping quality (Supplementary Figure 18b). First the probability  $\epsilon$  that the fragment was not from any of the HST in the cluster is calculated using the mapping quality  $q$ :

$$\epsilon = \max(\epsilon_{\min}, 10^{-q/10}), \quad (3)$$

where  $\epsilon_{\min}$  is the minimum noise probability. The motivation behind having a minimum is that mapping qualities are generally less reliable at higher values. The minimum noise probability is  $10^{-4}$  for all fragments except those that were labeled as low scoring, for which it is 1. Let  $A$  be the set of alignment paths (that is alignments) for a fragment. For

each alignment path  $a \in A$ , the likelihood of it being the correct path is calculated using its score  $s_a$  and length  $\ell_a$ :

$$L(a) = \psi_{\alpha} \left( \frac{\ell_a - \mu}{\sigma} \right) \exp(\lambda s_a), \quad (4)$$

where  $\lambda$  is a scaling factor that converts the alignment score into the log-likelihood of a pair hidden Markov model<sup>65</sup>,  $\psi_{\alpha}(x) = 2\phi(x)\Phi(\alpha x)$  is the density of a skew-normal distribution (with  $\phi$  and  $\Phi$  the density and distribution function of a standard normal distribution, respectively), and  $\mu$ ,  $\sigma$ , and  $\alpha$  are the location, scale, and shape parameters of the fragment length distribution modeled as a skew-normal distribution. For paired reads, these parameters are estimated from the alignment path lengths across all fragments that have (i) a mapping quality of at least 30, and (ii) the same length for all alignment paths. The fragment length distribution is omitted from the equation when the fragments are single-end reads. With this likelihood, we can compute the posterior probability that the fragment originated from a given HST. Let the set of all HST paths in the cluster be denoted by  $T$ , and let the set of HST paths an alignment path  $a$  is consistent with be denoted by  $T_a$ . The probability that the fragment (or alignment  $A$ ) originated from an HST is calculated as:

$$p_t = (1 - \epsilon) \cdot P(t|A) = (1 - \epsilon) \cdot \frac{P(A|t)P(t)}{\sum_{t \in T} P(A|t)P(t)} \quad (5)$$

with

$$P(A|t) \propto \max_{a \in A} \begin{cases} \frac{L(a)}{\bar{\ell}_t} & \text{if } t \in T_a \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

Here,  $\bar{\ell}_t$  is the effective transcript length for  $t$  calculated as  $\bar{\ell}_t = \ell_t - \mu_{\ell_t}$  (refs<sup>3,4</sup>). In turn,  $\mu_{\ell_t}$  is the mean of the fragment length distribution truncated to  $[1, \ell_t]$ , computed using a published formula<sup>66</sup>. The effective transcript length accounts for the fact that fragments cannot be sequenced from all positions owing to the size of the fragment. If the fragments are single-end reads, the fragment length distribution parameters used to calculate the effective length must be supplied by the user. The prior over HSTs  $P(t)$  is taken to be uniform. If the HST probability  $p_t$  is below  $10^{-8}$ , it is truncated to 0 to reduce storage.

We denote the set of all fragment probabilities in a cluster as  $F$  and the probabilities for a fragment  $i$  as  $F_i = (\epsilon, \mathbf{p})$ , where  $\mathbf{p}$  is the vector of probabilities over all  $T$  HSTs in the cluster. Many fragments will have very similar probabilities and can thus be collapsed to save computation resources and memory<sup>4</sup>. To do this we collapse two fragment probabilities  $F_i$  and  $F_j$  if they satisfy both of:

$$\begin{aligned} |\epsilon^i - \epsilon^j| &< 10^{-8} \\ |p_t^i - p_t^j| &< 10^{-8}, \quad \forall t \in T \end{aligned} \quad (7)$$

We also associate each set of collapsed fragments with  $c$ , the number of collapsed fragments in the set. The resulting set  $E$  of tuples  $(\epsilon, \mathbf{p}, c)$  is subsequently used to infer the expression of the HSTs in the pantranscriptome.

**Inferring haplotype-specific transcript expression.** RPVG quantifies the expression of the HSTs in the pantranscriptome using a nested inference scheme (Supplementary Figure 18c). This is done independently for each cluster. First, the distribution over haplotype combinations (that is, diplotypes) is inferred. The most probable haplotype combinations are then selected from this distribution and expression is inferred conditioned on the haplotypes. In the following, we will assume the



sample is diploid, but the equations and algorithms generalize to any ploidy.

The marginal distribution over diplotypes is approximated by assuming the haplotypes are identical for all transcripts in a cluster. The motivation behind this approximation is that most clusters cover only a small region (typically a gene) of the genome. However, this approximation can break down when there are partial haplotypes or recombination events in the cluster. Using the transcript and haplotype origin table provided by VG RNA, the HSTs in the cluster are first grouped by their haplotype origin. Note that since an HST can be consistent with more than one haplotype it can also belong to multiple groups. Next, groups with the same set of HSTs are collapsed, resulting in a set of unique haplotype groups.

Now let us denote the set of haplotype groups as  $H$ , with each group  $h \in H$  consisting of a set of HSTs. The objective is to infer the distribution over diplotypes  $d = \{h_1, h_2\}$  conditioned on the set of collapsed fragment probabilities  $E$ . The probability of a diplotype is defined as:

$$P(d|E) = P(\{h_1, h_2\}|E) \propto P(h_1)P(h_2) \prod_{(\epsilon, \mathbf{p}, \mathbf{c}) \in E} \left( \epsilon + \frac{1-\epsilon}{2} (P(\mathbf{p}|h_1) + P(\mathbf{p}|h_2)) \right)^c \quad (8)$$

and

$$P(\mathbf{p}|h) = \frac{\frac{1}{n} \sum_{t \in h} p_t}{\sum_{k \in H} \frac{1}{n} \sum_{t \in k} p_t} \propto \sum_{t \in h} p_t \quad (9)$$

where the prior probability of each haplotype group  $P(h)$  is proportional to the number of haplotypes in the group, and  $n$  is the number of transcripts in the cluster (the factors of  $\frac{1}{n}$  and  $\frac{1}{2}$  amount to an approximation that expression is uniform across all transcripts and the two haplotypes, respectively). This model is inspired by a similar haplotyping model used in DINDEL<sup>67</sup>.

The distribution over diplotypes is inferred by calculating  $P(d|E)$  for all pairs of haplotype groups  $h \in H$ . To reduce the space of haplotype combinations that need to be evaluated, RPVG uses a branch-and-bound-like algorithm, where diplotypes containing an improbable haplotype group are not evaluated. Instead, the probability of all diplotypes containing an improbable haplotype group is set to 0. A haplotype group  $h$  is labeled to be improbable if its optimal diplotype probability  $P(\{h, h_o\}|E)$  is  $s$  times lower than the current highest evaluated probability, where  $s$  is the minimum diplotype posterior probability threshold used in the next step in the inference. The optimal diplotype probability is defined as

$$P(\{h, h_o\}|E) \propto P(h) \prod_{(\epsilon, \mathbf{p}, \mathbf{c}) \in E} \left( \epsilon + \frac{1-\epsilon}{2} \left( P(\mathbf{p}|h) + \max_{h_o \in H} (P(\mathbf{p}|h_o)) \right) \right)^c \quad (10)$$

This value serves as an upper bound on the probability of any diplotype containing  $h$ .

The expression of the HSTs in the cluster is estimated using the inferred distribution over diplotypes. First, the set of diplotypes with a posterior probability of at least  $s = 10^{-3}$  is selected from the distribution  $P(d|E)$ . HST expression is inferred for each of the diplotypes in this set.

The following is repeated for each diplotype in the set. First, all HSTs that are consistent with at least one of the haplotypes in the diplotype are collected. We denote this HST subset  $T_s \subseteq T$  and define the likelihood over the relative expression values  $\alpha$  as

$$L(\alpha) = \prod_{(\epsilon, \mathbf{p}, \mathbf{c}) \in E} \left( \alpha_0 \epsilon + \sum_{t \in T_s} \alpha_t p_t \right)^c, \quad (11)$$

where  $\alpha_0$  is the expression value of an artificial ‘noise transcript’ that accounts for the possibility of mismapping. An expectation

maximization algorithm is used to find the (local) maximum likelihood estimate of the expression values. The algorithm iterates between assigning fractional fragment counts to the HSTs and the noise transcript, and updating the expression values. This is a well known algorithm that is used by many other transcript quantification tools<sup>1,3,4</sup>. The expression values are initialized uniformly and the expectation maximization algorithm is run until convergence or for a maximum of 10,000 iterations. The algorithm is considered converged if

$$\frac{|\alpha^i - \alpha^{i-1}|}{\alpha^i} \leq 0.001, \quad \forall \alpha \in \alpha : \alpha \geq 10^{-8}, \quad (12)$$

for ten consecutive iterations, where  $i$  is the index of the current iteration. This criteria is inspired by the one used by KALLISTO<sup>3</sup> and SALMON<sup>4</sup>. For the final maximum likelihood estimate, we truncate all the relative expression values below  $10^{-8}$  to 0.

After the expectation maximization step, RPVG can optionally run a Gibbs sampling step to quantify the uncertainty in the expression estimates. The Gibbs algorithm iteratively samples the assignment of each fragment to a HST (or the noise transcript), and the expression values  $\alpha$ , which are given a symmetric uniform Dirichlet prior with a concentration parameter of one. A similar algorithm is described in Li et al.<sup>1</sup> and Patro et al.<sup>4</sup>. First, 1,000 diplotypes are sampled from  $P(d|E)$ , with a Gibbs sampler being run for each unique sampled diplotype. Each sampler is initialized on the maximum likelihood estimate from the expectation maximization algorithm and the number of samples of expression values collected is equal to the number of times the diplotype was sampled. This results in a total of a 1,000 collected Gibbs samples of expression values for each cluster. In addition, we thin each Gibbs chain and only collects a sample of expression values at every 25th Gibbs iteration. This is done to reduce autocorrelation between samples.

RPVG provides both a joint and marginal output of the inferred probabilities and expression values. The joint output contains the inferred posterior probabilities over HST combinations (that is, diplotypes) and their corresponding estimated expression values. Only combinations with a probability of at least  $10^{-8}$  are written to this output. The marginal output contains the haplotype probability and estimated expression value for each HST in the pantranscriptome. The haplotype probability is calculated as the sum of posterior probabilities over all diplotypes that include the HST. The HST expression is similarly calculated as the sum of the estimated expression values over all diplotypes that include the HST, weighted by their posterior probability. The expression of the noise transcript is aggregated across all clusters into a single artificial transcript called ‘Unknown’.

### Transcript quantification evaluation

We compared the quantification accuracy of RPVG against three other transcript quantification tools: KALLISTO; SALMON; and RSEM. Haplotype-specific transcript indexes for KALLISTO, SALMON, and RSEM were built from the HST sequence FASTA files generated by VG RNA. For the real data, the 104 full-length mitochondrial and scaffold transcripts in the GENCODE v29 annotation were added to the pantranscriptomes. SALMON indexing was run with duplicates kept and, on the real data, the reference genome was given as a decoy. The Bowtie2 mapper was used in RSEM with the maximum number of alignments per read increased to 1,000. The transcript expression was estimated using default parameters for all methods, except for the real data where strand-specific inference was enabled. KALLISTO and SALMON were run without bias correction. RSEM was only run on the NA12878 personal-sample-specific transcriptome and the ‘Europe (excl. CEU)’ pantranscriptome, as it did not scale to the two largest pantranscriptomes.

RPVG was run using default parameters and with three different types of alignments inputs: the standard multipath alignments from VG MPMAP and single-path alignments from VG MAP and VG MPMAP.



The VG MPMAP single-path alignments were generated by finding the best scoring path in the multipath alignments using VG VIEW. The fragment length distribution parameters estimated by VG MPMAP were given as input to RPVG when using the VG MAP alignments. RPVG was run with a ploidy of two for all read sets, including CHM13. All HSTs with a haplotype probability below 0.8 were filtered from the RPVG output.

The main expression results were obtained using the ENCSR000AED, replicate 1 data (see Supplementary Table 4): Fig. 3; Extended Data Figs 7 and 9; and Supplementary Figs 8–13. The SRR1153470 and CHM13 read data (see Supplementary Table 4) was used to optimize the parameters of RPVG: Supplementary Figs 14 and 15.

For the ENCSR000AED and SRR1153470 data, which are both NA12878 cell lines, we compared the quantified HSTs to the haplotypes of NA12878 from the 1000GP data. We considered an HST consistent with these haplotypes if it matched the sequence of one of the two possible NA12878 haplotype versions of the transcript. Biopython was used to parse and compare HST FASTA sequences<sup>68</sup>. The haplotyping performance of each method was then estimated by comparing the number and fraction of quantified HSTs with positive expression that were consistent.

We used TPM to measure expression. For the simulated data we recalculated the TPM value for all methods to ensure that there was no bias towards RSEM, which was used to estimate the expression profile employed by VG SIM to parameterize the HST expression values. The TPM value depends on the effective transcript length, which is not calculated in the same manner for each method. Therefore, if this is not corrected, methods that estimate the effective transcript length more similarly to RSEM will have an advantage that does not depend on their ability to predict correct expression values. The true fragment length distribution parameters and the effective transcript length approach employed by RPVG (similar to KALLISTO and SALMON) was used when recalculating the TPM values.

The ability of the method to predict the correct expression value was evaluated using the simulated data for which the true expression is known. The true expression values were calculated from a table provided by VG SIM, which indicates the transcript of origin for each read. The simulated TPM values were calculated in the same manner as described above. We used both Spearman correlation and mean absolute relative difference (MARD) to quantify concordance between estimated and true expression.

The CHM13 cell line is effectively haploid, so only a single HST is expected to exist for each transcript. We used this feature of the data to measure the haplotype inference performance of each method on the T2T CHM13 data. We defined each HST as either major or minor. Major HSTs were defined as the highest-expressed haplotype for each transcript; the rest were defined as minor. The fraction of expression from minor HSTs is a lower bound on the fraction of incorrectly inferred transcript expression. Accordingly, we used the number of major and minor transcripts that each method predicted to be expressed to compare their haplotype inference performance.

To evaluate allele-specific expression (ASE) estimation, we converted the simulated and estimated HST expression values to allele-specific read counts for the NA12878 variants. These were calculated by dividing the expression values with the corresponding transcript length and multiplying by twice the read length (to account for paired-end sequencing). In addition, we inferred allele-specific read counts for the same NA12878 variants using the WASP<sup>19</sup> pipeline with STAR alignments. Using both simulated and inferred allele-specific read counts, we next labeled heterozygotic variants with at least one read in the simulated data as showing significant ASE using a two-sided binomial test with *P* values adjusted using the Benjamini–Hochberg procedure and an FDR  $\alpha = 0.1$ . We took the hypothesis tests of the true simulated read counts to be the truth labels. The sets of labeled heterozygotic variants were lastly compared between simulation and methods to produce ASE true and false positive rates.

We assessed the robustness of RPVG to admixture on real data using an indirect proxy. We applied the MPMAP–RPVG pipeline to two samples from the same study<sup>30</sup>: one of European American ancestry (SRR12765534) and the other of African American ancestry (SRR12765650). We expect the African American individual to be more highly admixed. We then measured the proportion of marginal posterior expression assigned to the two most highly expressed HSTs by summing over diplotypes. This is a lower bound on error, since for the majority of genes without a copy number alteration, there can only be two copies of the gene. We then computed the proportion of transcripts for which the two highest-expressed HSTs accounted for at least a given threshold proportion of the total expression. We repeated this analysis for different threshold values and stratified the results by minimum transcript expression.

### HLA pantranscriptome construction and typing

We evaluated the ability of the VG MPMAP–RPVG pipeline to type HLA alleles. To start, we constructed a set of HLA haplotypes using gene allele sequences from the IPD-IMGT/HLA database (release 3.43.0)<sup>32</sup>. Many of the alleles in the database are partial and do not cover the corresponding entire gene, with a large fraction of them only covering the coding sequence or just the antigen recognition site (ARS) exons. Since haplotypes covering whole genes including introns are needed to construct a pantranscriptome using the VG RNA pipeline, we first imputed the missing coding sequence with the closest complete allele for all the partial alleles using HLASEQLIB<sup>23</sup>.

Next we used the reference to extend the alleles into full haplotypes. We padded each allele with 10,000 reference bases on both sides using the corresponding genes coding start and end location in the GENCODE v29 transcript annotation to ensure that the allele sequences would align to the correct genes. The padded HLA alleles were then aligned to a spliced pangenome graph using VG MPMAP in long read, single-path mode. The resulting alignments were projected to the reference genome using VG SURJECT and used to determine the location of splice junctions in the allele sequences. The reference sequences of the corresponding introns were added to the allele to produce haplotype sequences covering the whole gene. The intron sequences were only added for junctions that were within 2 bases of an annotated splice junction to ensure that genomic deletions were not mistakenly interpreted as splice junctions.

These haplotypes were then used to create HLA pantranscriptomes. First, the haplotypes were mapped against the same spliced variation graph using VG MPMAP in long read, single-path mode. The resulting alignments were used to update the graph with the variation in the haplotypes using VG AUGMENT. Using these haplotypes and VG RNA, we created two HLA pantranscriptomes (see Supplementary Table 3): “HLA (main)” consisting of five of the main and most variable HLA genes and “HLA (10)” consisting of ten HLA genes of which all had at least a 100 haplotypes. Null alleles that have been shown not to be expressed were not included in the construction of the pantranscriptomes. In addition, transcript *HLA-B-258* was also not used as it was covering both *HLA-B* and *HLA-C*. This transcript have been removed in later versions of the GENCODE annotation. “HLA (main)” was built using the “1000GP (all, excl. CEU)” spliced variation graph, whereas “HLA (10)” was built using the “1000GP (all)” graph.

Using the “HLA (main)” pantranscriptome, we first optimized the default parameters of RPVG for HLA typing using six RNA-seq samples from the Geuvadis dataset: NA07051, NA11832, NA11840, NA11930, NA12287, and NA12775<sup>34</sup> (see Supplementary Table 4). All samples were from the CEU population, which was excluded in the variation graph construction. We compared the inferred alleles to the typing results from Gourraud et al.<sup>35</sup> and Abi-Rached et al.<sup>36</sup>, both of which are available on the 1000GP homepage (<https://www.internationalgenome.org/category/hla/>). Similar to the quantification evaluation, HSTs with a haplotype probability below 0.8 were filtered before evaluation.

An expressed HST was regarded as correct if its corresponding HLA allele matched one of the two studies. When evaluating diplotyping performance both HLA alleles needed to be correct and also match both ground-truth alleles in the same study. To improve accuracy and speed, we adjusted the maximum number of standard deviations from the mean allowed for the fragment length, the maximum allowed score difference to the best alignment and the threshold for filtering alignments compared to the optimal score.

Next, we ran RPVG using the optimized parameters on two different sets of RNA-seq data: ten randomly selected CEU samples from Geuvadis that were not used in the optimization and three parent–child trios from the 1000GP sequenced as part of the Human Genome Structural Variation Consortium (HGSVC)<sup>33</sup> (see Supplementary Table 4). The Geuvadis and HGSVC datasets were run on the “HLA (main)” and “HLA (10)” pantranscriptomes, respectively. For the Geuvadis data, we used the same two studies as described above to determine typing accuracy, whereas for the HGSVC data trio concordance was used. For the Geuvadis datasets, typing accuracy for three different levels of HLA resolution were evaluated: one field, two field and G groups. G groups are defined as alleles that have identical nucleotide sequences across the ARS exons and were used to distinguish ambiguous alleles in the Gourraud et al. study.

### Variant genotyping and effect prediction

We used RNA-seq data from five randomly selected tissues from the same individual to estimate allele concordance across datasets. We also demonstrate the ability to genotype variants with potential effects on functional elements. All five datasets are available from the ENCODE project<sup>28,29</sup> (see Supplementary Table 4). For each RNA-seq dataset, all technical replicates were combined and the VG MPMAP–RPVG pipeline was run using the “Whole” pantranscriptome (see Supplementary Table 3). The pipeline was run with default parameters except for RPVG, for which it was specified that the RNA-seq data was strand specific.

The RPVG HST expression estimates were converted to variant allele-based expression values for the downstream analyses. To do this, all exonic variants were first annotated with the transcripts overlapping them using BCFTOOLS. These annotated variants were then used together with the original haplotypes to translate each HST to its corresponding set of variant alleles. Using this translation, we computed the expression of each variant allele as the sum of the expression over HSTs that contained the allele. The haplotype probability values were similarly computed as a sum over the HSTs that contained the allele. However, for diplotypes where the alleles were called homozygous, only the probability of one of the haplotypes was added. This ensured that the corresponding alleles were only counted once per diplotype sample similar to the haplotype probabilities.

Using these results for each of the five tissues, we estimated the number of expressed variant alleles and the allele concordance between the tissues. For this analysis, we filtered all alleles with a probability below 0.8. An allele was defined to be expressed if it had non-zero expression in at least one tissue, and a variant was defined as expressed if at least one of its alleles, including the reference, had non-zero expression. Next, we estimated the consistency in whether an allele was expressed or not between tissues. To account for alternative expression and splicing across tissues, we only considered tissues for which the corresponding variant was expressed. An allele was then said to be concordant across tissues if it was either consistently expressed or consistently not expressed across all tissues that had the corresponding variant expressed (see Extended Data Figure 10). Variants that were only expressed in a single tissue were excluded for the concordance estimation. Since both alleles might not have been sequenced by chance for lowly expressed exons, we repeated the analyses for different thresholds of variant expression. Finally, the homopolymer length of each variant was calculated by counting the maximum number of consecutive identical bases in each direction from the variant start site.

Next, we predicted the effect of the expressed variants on functional elements, such as transcript and protein sequences. This was done using the Ensembl Variant Effect Predictor (VEP) toolset<sup>38</sup>. VEP was run on all expressed variant alleles with a probability of at least 0.8 and a variant expression value of at least 5 TPM. The reason for the latter threshold was that we wanted to show the effects for exons with a non-negligible expression and also minimize genotyping error from allelic dropout. The GENCODE v29 transcript annotation was used for VEP with conversion to minimal variant representation enabled. Predicted consequences that did not have an impact rating of moderate or high were filtered. For variants with both moderate and high impact consequences, only the high impact ones were shown.

### Demonstration of analyzing genomic imprinting

We obtained RNA-seq datasets from samples NA11832, NA11930, NA12775, and NA12889 from the Geuvadis dataset<sup>34</sup> and ran them through the VG MPMAP–RPVG pipeline (see Supplementary Table 4). Each sample had two accessions, which were combined into one dataset. We also obtained and analyzed data from sample NA12878 from the ENCODE Project (ENCSR000AED replicate 1)<sup>28,29</sup>. These samples are all unrelated. All parameters used were identical to those used in the real data evaluations of VG MPMAP and the only difference for RPVG was that Gibbs sampling was enabled. The Geuvadis samples were used to troubleshoot the analysis and identify potentially interesting genes to highlight in the demonstration. The analyses were then repeated on the ENCODE sample. This design reduces the risk of identifying noise as signal. Only the results of the final analysis are the ones reported in the Results section.

To confirm that the pipeline could detect previously known ASE, we looked for signatures of imprinting in the 20 genes with the most statistically significant parent-of-origin ASE in the study by Zink et al.<sup>17</sup> (from Supplementary Table 6). One of these genes, *RP11-69E11.4*, had since been removed from the GENCODE database, so we excluded it from the analysis. The Zink et al. study analyzed ASE on individual SNVs. To make our results comparable to theirs, we translated the RPVG HST-based expression quantification into a corresponding variant allele-based expression quantification using the same approach as described above. The expression of each allele was computed as the sum of the expression of each HST that contained the allele.

We decided to highlight the haplotype-specific expression of the *NAA60* gene in depth because it consistently showed monoallelic expression for both haplotypes across different isoforms in the initial exploratory datasets. To identify the haplotype of origin for different HSTs, we compared the variants associated with each HST (using the table from VG RNA) to the haplotypes of the sample from the 1000GP dataset. Equal-tailed credible intervals were approximated using RPVG's Gibbs sampling method.

### Reporting summary

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

### Data availability

All data used in this study are available at <https://github.com/jonas-sibbesen/vgrna-project-paper>. Data that are available from public repositories are provided as web links only. Accession numbers are included when relevant, and accession numbers for sequencing data are also listed in Supplementary Table 4. The repository also includes all spliced pangenome graphs and pantranscriptome haplotype-specific transcript sets, which may be freely used in other projects. Mapping benchmark tables and haplotype-specific expression estimates are archived in Zenodo (<https://doi.org/10.5281/zenodo.7234454>).

## Code availability

The source code for VG and RPKG is publicly available at <https://github.com/vgteam/vg> and <https://github.com/jonassibbesen/rpkg>, respectively. Both tools are licensed under the MIT License. A full list of the versions of all computational tools used is available in Supplementary Table 6. All bash scripts with exact command-lines used to generate the results are available at <https://github.com/jonassibbesen/vgrna-project-paper>. This repository also includes the custom C++, Python, and R scripts used for analysis and plotting, together with references to Docker containers and log files from the analyses.

## References

53. Danecek, P. et al. Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008 (2021).
54. Eizenga, J. M. et al. Efficient dynamic variation graphs. *Bioinformatics* **36**, 5139–5144 (2020).
55. Gagie, T., Navarro, G. & Prezza, N. Fully functional suffix trees and optimal text searching in BWT-runs bounded space. *J. ACM* **67**, 1–54 (2020).
56. Sirén, J. Indexing variation graphs. In *2017 Proc. 19th Workshop on Algorithm Engineering and Experiments (ALENEX)* 13–27 (SIAM, 2017).
57. Chang, X., Eizenga, J., Novak, A. M., Sirén, J. & Paten, B. Distance indexing and seed clustering in sequence graphs. *Bioinformatics* **36**, 146–153 (2020).
58. Paten, B. et al. Superbubbles, ultrabubbles, and cacti. *J. Comput. Biol.* **25**, 649–663 (2018).
59. Eades, P., Lin, X. & Smyth, W. F. A fast and effective heuristic for the feedback arc set problem. *Inf. Process. Lett.* **47**, 319–323 (1993).
60. Lee, C., Grasso, C. & Sharlow, M. F. Multiple sequence alignment using partial order graphs. *Bioinformatics* **18**, 452–464 (2002).
61. Burset, M., Seledtsov, I. A. & Solovyev, V. V. Analysis of canonical and non-canonical splice sites in mammalian genomes. *Nucleic Acids Res.* **28**, 4364–4375 (2017).
62. Li, H. et al. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
63. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
64. Wala, J. & Beroukhi, R. SeqLib: a C++ API for rapid BAM manipulation, sequence alignment and sequence assembly. *Bioinformatics* **33**, 751–753 (2016).
65. Karlin, S. & Altschul, S. F. Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc. Natl Acad. Sci. USA* **87**, 2264–2268 (1990).
66. Flecher, C., Allard, D. & Naveau, P. Truncated skew-normal distributions: moments, estimation by weighted moments and application to climatic data. *Metron* **68**, 331–345 (2010).
67. Albers, C. A. et al. Dindel: accurate indel calls from short-read data. *Genome Res.* **21**, 961–973 (2011).

68. Cock, P. J. A. et al. Biopython: freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).

## Acknowledgements

Research reported in this publication was supported by the National Human Genome Research Institute of the National Institutes of Health under award numbers U01HG010961, R01HG010485, U41HG010972, U24HG011853, and OT2OD026682 to B.P. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. The work of J.A.S. was supported by the Carlsberg Foundation. We thank the ENCODE Consortium, the Thomas Gingeras Laboratory (Cold Spring Harbor Laboratory), the Ali Mortazavi Laboratory (University of California Irvine) and the Joe Ecker Laboratory (Salk Institute for Biological Studies) for generating and sharing the ENCODE data used in this study. We would also like to thank M. Dennis (University of California Davis) for generating and providing access to the CHM13 RNA-seq data on behalf of the T2T consortium. Finally, we would like to thank J. Monlong and G. Hickey for feedback on the manuscript, and everybody else in the VG Team.

## Author contributions

J.A.S. and J.M.E. developed software, designed and carried out experiments, analyzed data, and wrote the paper. A.M.N., J.S., X.C., and E.G. contributed to developing the software. B.P. contributed to project conceptualization, supervised the research, and edited the paper.

## Competing interests

The authors declare no competing interests.

## Additional information

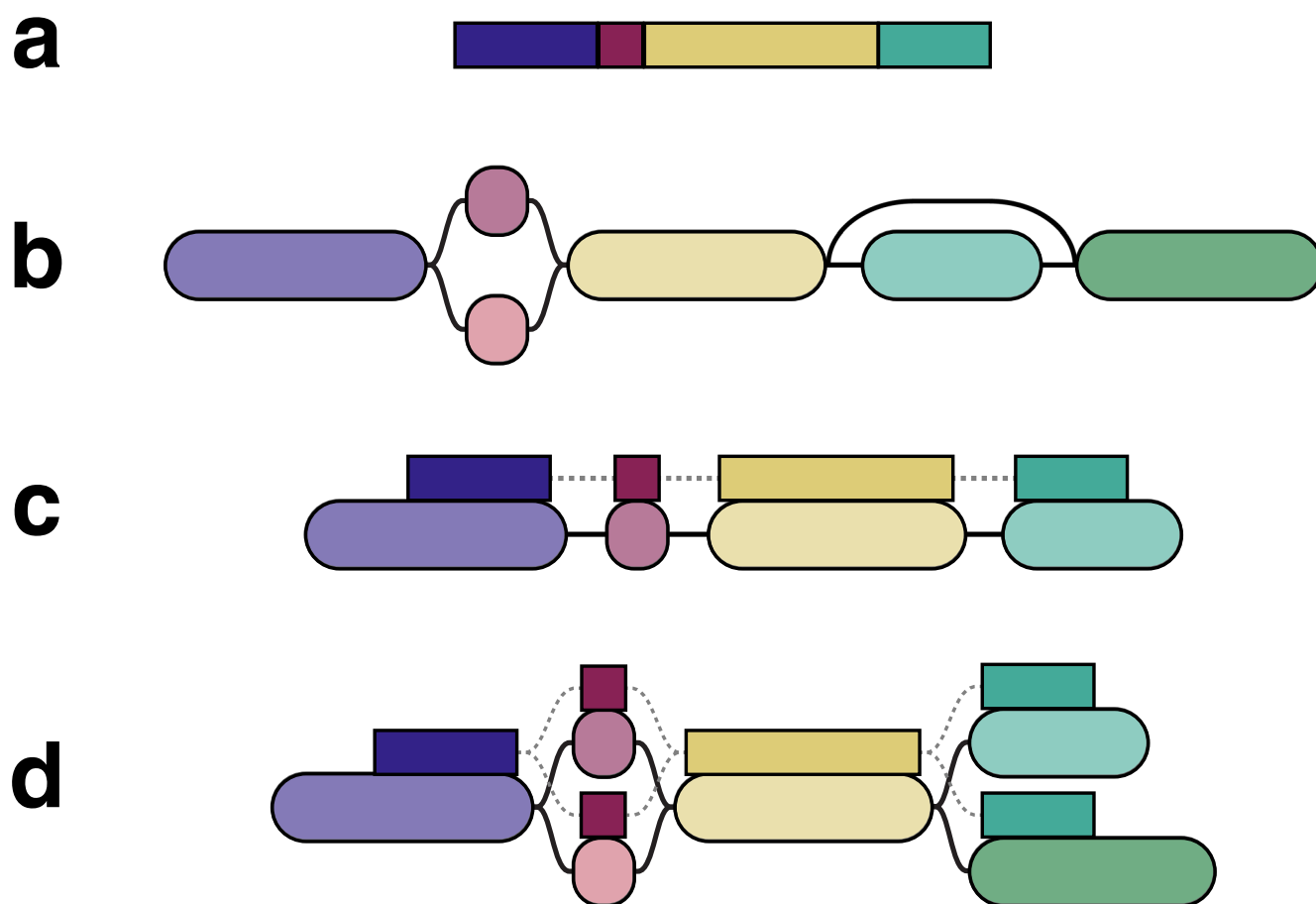
**Extended data** is available for this paper at <https://doi.org/10.1038/s41592-022-01731-9>.

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41592-022-01731-9>.

**Correspondence and requests for materials** should be addressed to Benedict Paten.

**Peer review information** *Nature Methods* thanks Michael Love, Harold Pimentel, and the other, anonymous, reviewer for their contribution to the peer review of this work. Primary Handling Editor: Lin Tang, in collaboration with the *Nature Methods* team.

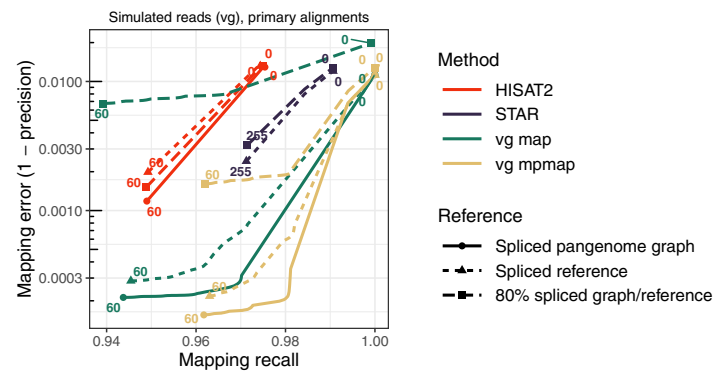
**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).



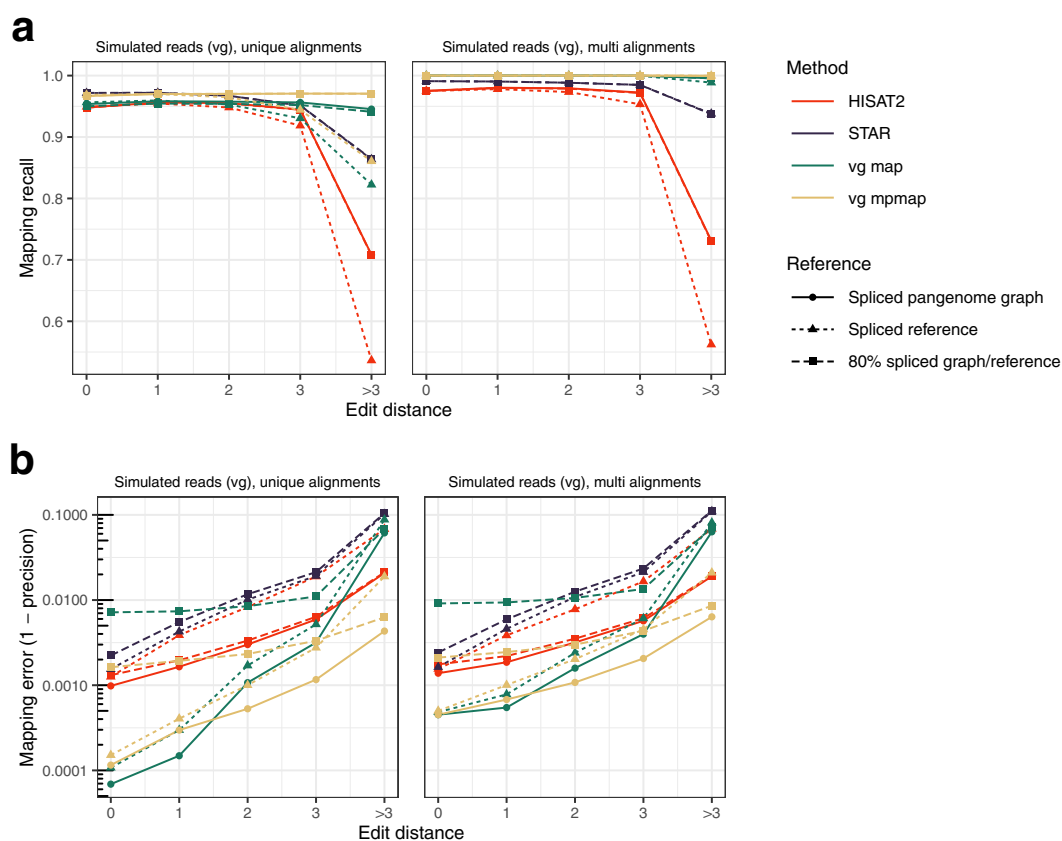
**Extended Data Fig. 1 | Diagram of a multipath alignment.** A diagrammatic comparison between the multipath alignment output of VG MPMAP and the single-path alignment output of other graph aligners (such as VG MAP). **a** A read and **b** a sequence graph, which have been colored to indicate which

parts of the read could plausibly align to which parts of the graph. **c** A single-path alignment. The read sequence is aligned to one path from the graph. **d** A multipath alignment. The alignment can split and rejoin to express the alignment uncertainty to different paths in the graph.



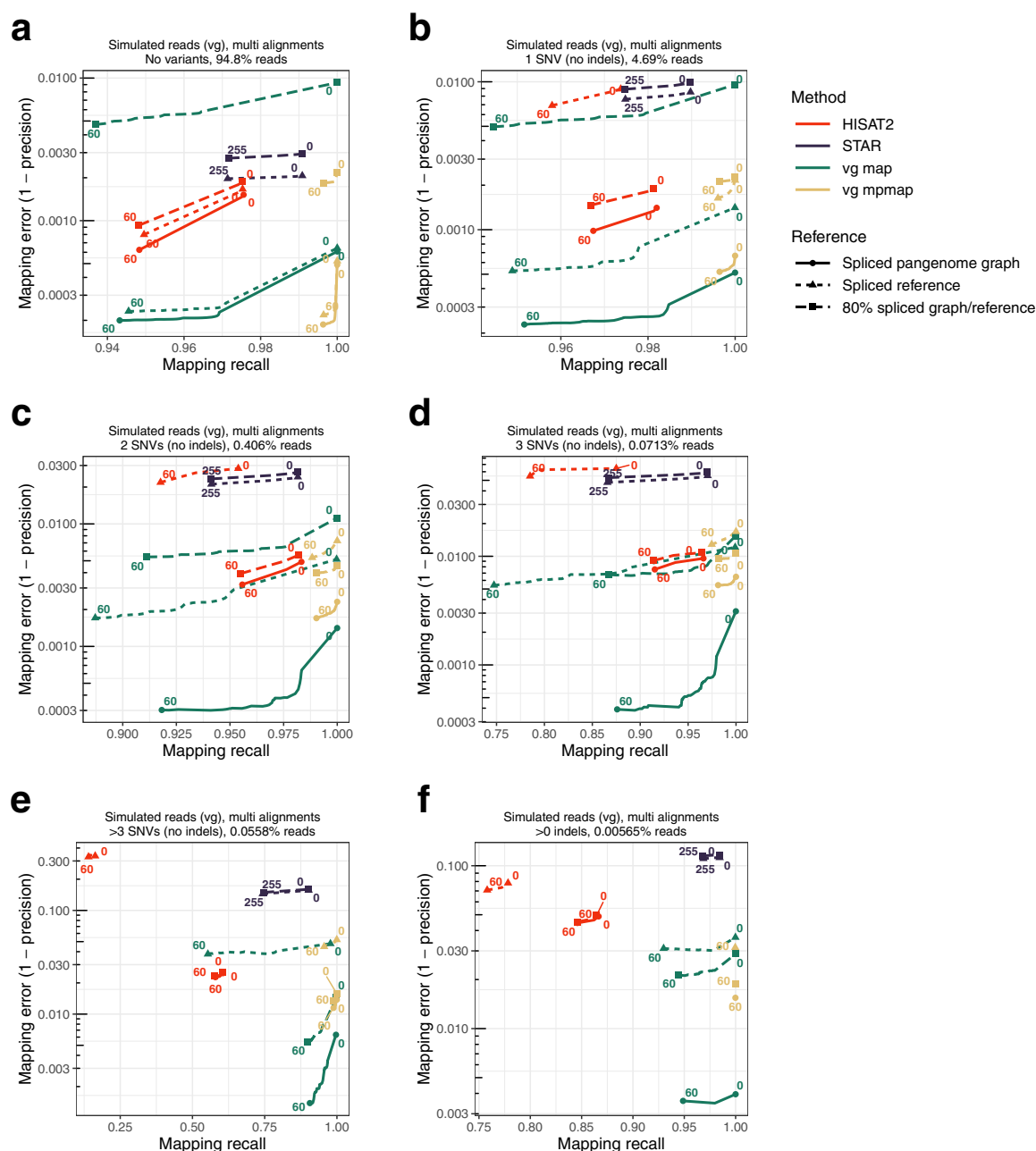


**Extended Data Fig. 2 | Mapping benchmark for primary alignments using RNA-seq data from NA12878.** Mapping error and recall for VG MPMAP and three other methods using simulated Illumina data. Colored numbers indicate different mapping quality thresholds. Reads are considered correctly mapped if their primary alignments cover 90% of the true reference sequence alignment.



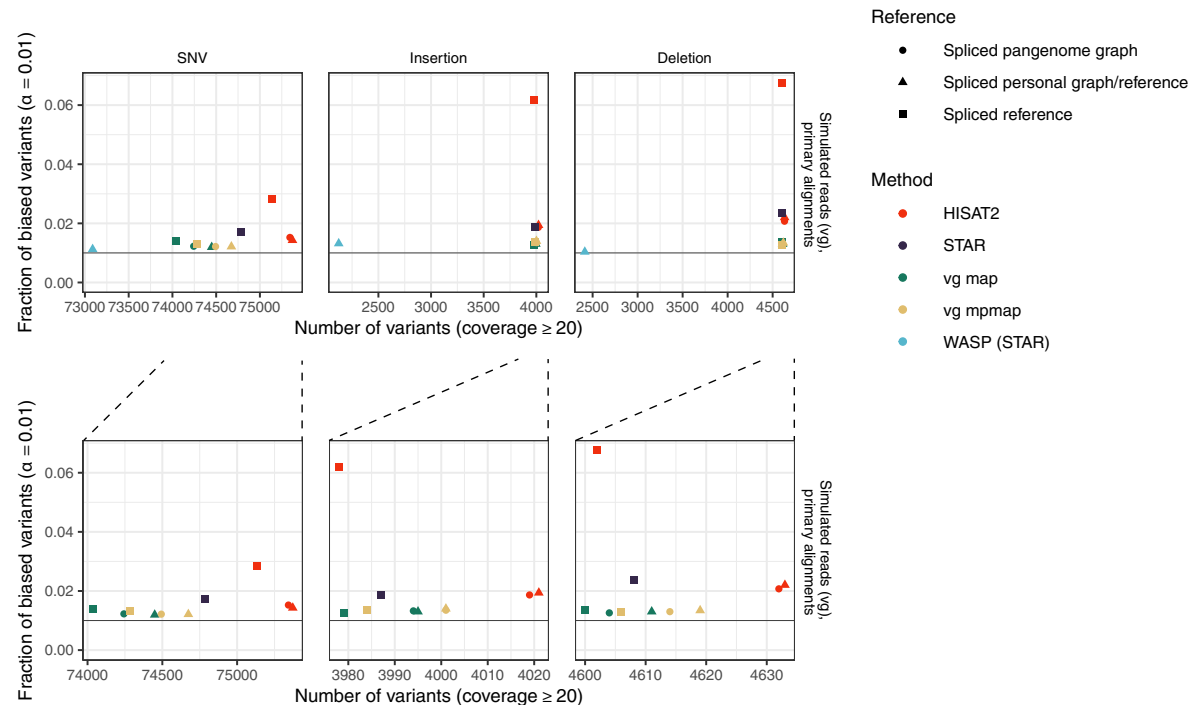
**Extended Data Fig. 3 | Mapping benchmark stratified by edit distance using RNA-seq data from NA12878.** Mapping recall (a) and error (b) for VG MPMAP and three other methods using simulated Illumina data as a function of edit

distance. Unique alignments are primary alignments with a mapping quality of at least 30. Reads are considered correctly mapped if their alignments cover 90% of the true reference sequence alignment.



**Extended Data Fig. 4 | Mapping benchmark stratified by non-reference variants using RNA-seq data from NA12878.** Mapping error and recall for VG MPMP and three other methods using simulated Illumina data. Colored numbers indicate different mapping quality thresholds. Reads are considered correctly mapped if one of their multi-alignments covers 90% of the true

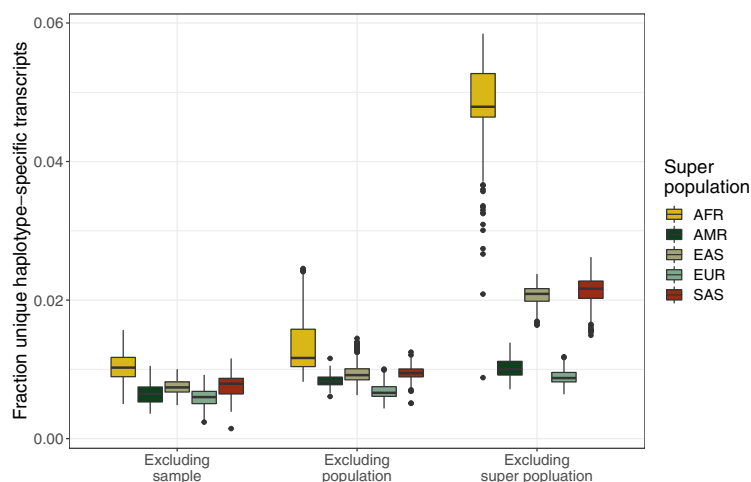
reference sequence alignment. Reads are stratified into those that **a** contain no variants, **b** contain no insertions or deletions (indels) and one single nucleotide variant (SNV), **c** contain no indels and two SNVs, **d** contain no indels and three SNVs, **e** contain no indels and more than three SNVs, and **f** contain any indels.



**Extended Data Fig. 5 | Allelic bias benchmark using RNA-seq data from NA12878.** Allelic mapping bias for VG MPMAP and four other methods using simulated Illumina RNA-seq reads, which were simulated without allelic bias. STAR was used as the aligner for the WASP pipeline. The WASP (STAR) pipeline were provided the 1000GP NA12878 haplotypes as input. The number of variant

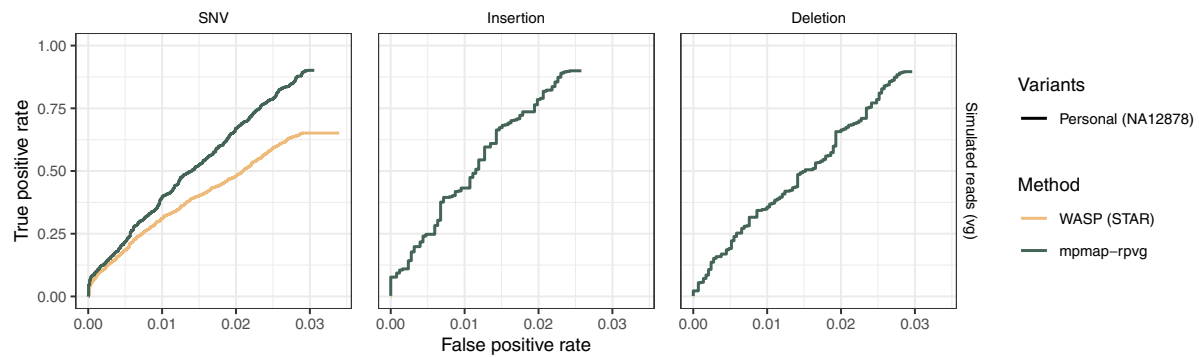
sites with coverage at least 20 is plotted against the observed rate of false positive hypothesis tests of allelic skew (two-sided binomial test,  $\alpha = 0.01$ ). Coverage was calculated from primary alignments with a mapping quality value of at least 30. The bottom row shows a zoomed view without WASP (STAR).





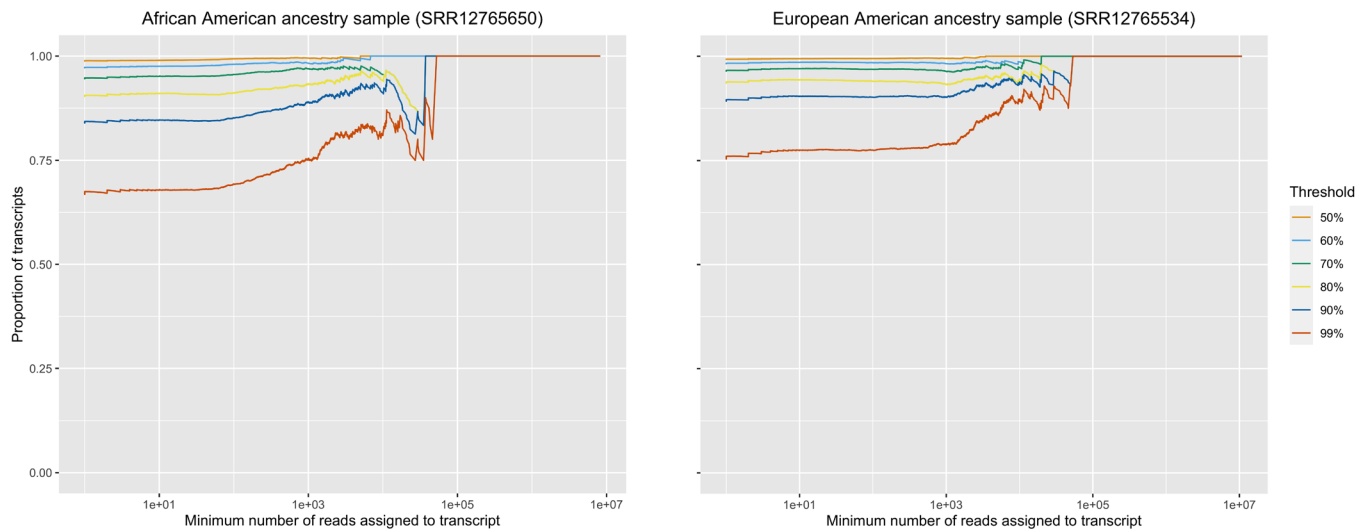
**Extended Data Fig. 6 | Haplotype-specific transcript uniqueness in a 1000 Genomes Project.** The fraction of HSTs that are unique to each of the 2504 samples in the 1000 Genomes Project (1000GP) when compared to different subsets of samples in the 1000GP. Left box plots show the fraction unique when comparing to all other samples, middle box plots show the fraction unique when comparing to all other samples excluding the samples' population, and right box plots show the fraction unique when comparing to all other samples excluding

the samples' super population. AFR: African ( $n = 661$ ), AMR: Admixed American ( $n = 347$ ), EAS: East Asian ( $n = 504$ ), EUR: European ( $n = 503$ ), SAS: South Asian ( $n = 489$ ). The horizontal line in the boxes corresponds to the median, and the box bounds (inter-quartile range) to the 25th and 75th percentile. The whiskers extend to the minimum and maximum value, but no further than 1.5 times the inter-quartile range from the box bounds. Values outside the whiskers are displayed as points.



**Extended Data Fig. 7 | Allele-specific expression benchmark using RNA-seq data from NA12878.** Allele-specific expression (ASE) results comparing the MPMAP-RPVG pipeline against WASP (with STAR as the aligner) using simulated data. Shows true positive rate and false positive rate of ASE significance for different thresholds of variant read count in the simulated data. Variants were defined as showing significant ASE using a two-sided binomial test of the allele-specific read counts with p-values adjusted using the Benjamini-Hochberg

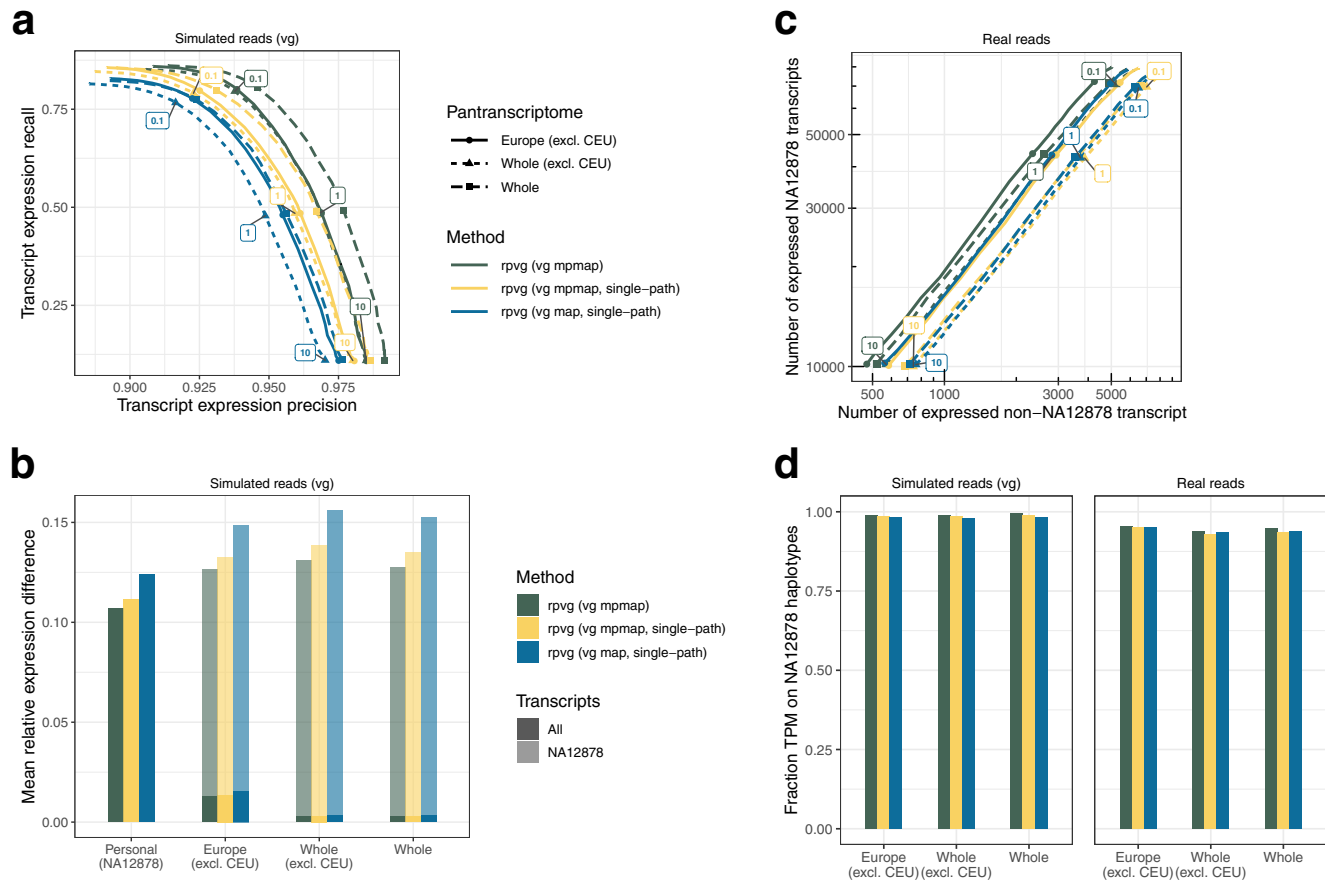
procedure and a False Discovery Rate (FDR)  $\alpha = 0.1$ . All heterozygotic NA12878 variants from the 1000 Genomes Project (1000GP) with at least one read in the simulated data were used for the benchmark. For the MPMAP-RPVG pipeline, we used the personal transcriptome generated from the 1000GP NA12878 haplotypes (Supplementary Table 3). WASP was provided the 1000GP NA12878 haplotypes as input. Note, we only used WASP for bias correction and allele-specific read counting, and not its downstream inference method.



**Extended Data Fig. 8 | Proportion of marginal expression attributed to  $\leq 2$  HSTs of a transcript.** For an African American individual (left) and a European American individual (right), the proportion of transcripts for which the marginal expression has at least  $X$  proportion assigned to  $\leq 2$  HSTs is shown for various

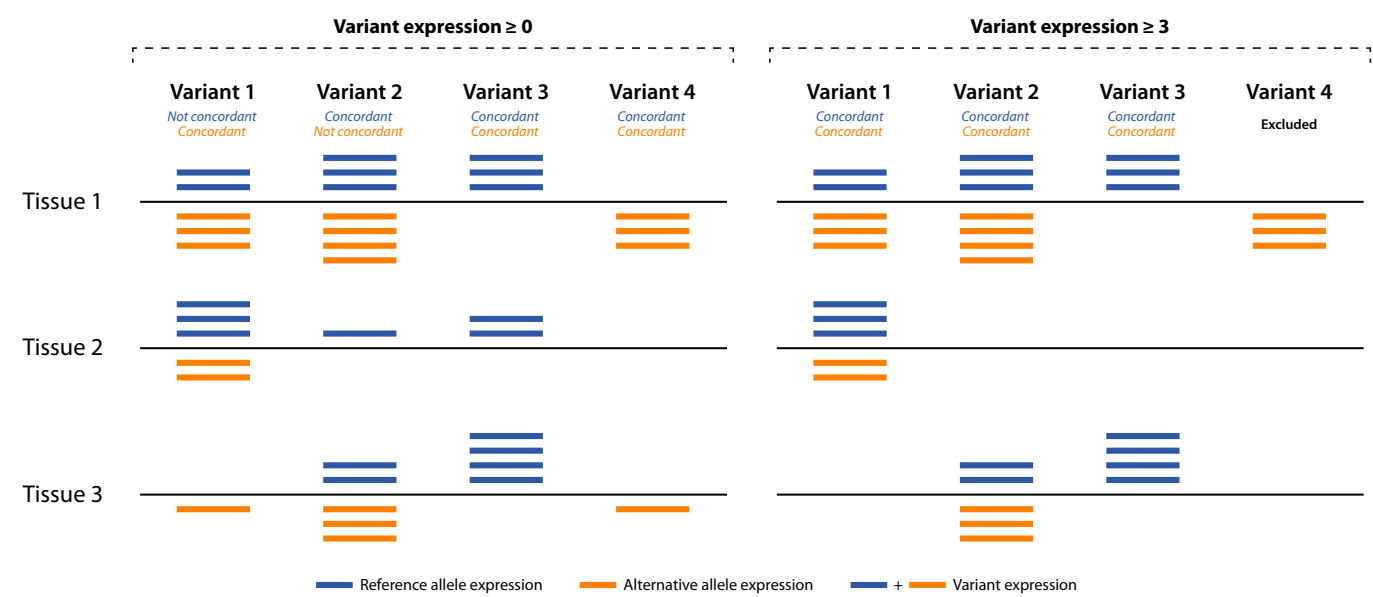
values of  $X$ . Colors correspond to different thresholds on the proportion of marginal expression. A pantranscriptome generated from all 1000 Genomes Project haplotypes were used for the evaluation ("Whole" in Supplementary Table 3). Transcripts with fewer than 1 inferred read are omitted.





**Extended Data Fig. 9 | Multipath alignment benchmark using RNA-seq data from NA12878.** Haplotype-specific transcript (HST) quantification results comparing RPVG with single-path and multipath alignments from VG MPMPAP and VG MAP as input using simulated and real Illumina data. For details on the pantranscriptomes used see Supplementary Table 3. The VG MPMPAP single-path alignments were created by finding the best scoring path in each multipath alignment. **a** Recall and precision of whether a transcript is correctly assigned nonzero expression for different expression value thresholds (colored numbers for “Whole (excl. CEU)” pantranscriptome) using simulated data. Expression is

measured in transcripts per million (TPM). **b** Mean absolute relative expression difference (MARD) between simulated and estimated expression (in TPM) for different pantranscriptomes using simulated data. MARD was calculated using either all HSTs in the pantranscriptome (solid bars) or using only the NA12878 HSTs (shaded bars). **c** Number of expressed transcripts from NA12878 haplotypes shown against the number from non-NA12878 haplotypes for different expression value thresholds (colored numbers) using real data. **d** Fraction of transcript expression (in TPM) assigned to NA12878 haplotypes for different pantranscriptomes using simulated (left) and real (right) data.



**Extended Data Fig. 10 | Examples of allele expression concordance across tissues.** A set of examples showing allele concordance across tissues using two different variant expression thresholds. Only three tissues are used in the example for simplicity. Blue and orange bars correspond to reference and alternative allele expression, respectively. Variant expression is calculated as the sum of the two alleles. An allele is defined as concordant if it is either consistently expressed or consistently not expressed across all tissues for which

the corresponding variant is expressed. Using this definition all alternative alleles except for the allele in variant 2 are defined as concordant when the minimum variant expression threshold is set to 0. If the variant expression threshold is increased to 3, the alternative allele in variant 2 becomes concordant since tissue 2 will be filtered for this variant. Moreover, variant 4 will be excluded due to tissue 3 being filtered since at least two expressed tissues are needed to compute concordance.

## 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 our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- | n/a                                 | Confirmed  |
|-------------------------------------|--|
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided<br><i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i>   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A description of all covariates tested   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A full description of the statistical parameters 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) |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. $F$ , $t$ , $r$ ) with confidence intervals, effect sizes, degrees of freedom and $P$ value noted<br><i>Give <math>P</math> values as exact values whenever suitable.</i>                            |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's $d$ , Pearson's $r$ ), indicating how they were calculated   |

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

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

Data collection

Data analysis   
  
bcftools v1.9 & v1.11  
samtools v1.9  
bedtools v2.29.1  
seqtk v1.3  
HISAT2 v2.2.0 & v2.2.1  
STAR v2.7.9a  
WASP v0.3.4  
Bowtie2 v2.3.5.1 & v2.4.4  
RSEM v1.3.1 & v1.3.3  
Kallisto v0.46.1 & v0.46.2  
Salmon v1.2.1 & v1.5.2  
VEP v103.1  
SeqLib 08771285  
Biopython v1.71 & v1.77  
hlaseq v0.0.2  
rpv 1d91a9e3



vgna-project-scripts 71442ea4, 94176204, 34d563ba, 7eb08153, 93bc0a90, 2220bb08, f9db749e, 9a079ecf & b08defd4

vg construct, vg convert, vg rna, vg ids, vg index, vg stats, & vg gbwt (Constructing graphs and pantranscriptomes) v1.23.0, c861e23e, 8ff022c3 & c4bbd63b  
 vg gbwt (Constructing GBWT r-index) 883f0f87 & c4bbd63b  
 vg snarls, vg prune, & vg index (Constructing distance and GCSA graph index) 8ff022c3 & c4bbd63b  
 vg mpmap, vg surject, vg stats, & vg augment (Augmenting graph with HLA haplotypes) c4bbd63b  
 vg paths & vg surject & (Creating reference transcript alignments for mapping benchmark) c861e23e  
 vg view & vg sim (Simulating reads from transcript paths) 765d2215  
 vg map & vg mpmap (Mapping reads to graph) 2cea1e25  
 vg inject, vg view, vg paths, & vg surject (Converting alignments between graph (GAM) and reference (BAM)) 385fd636 & 2cea1e25  
 vg stats, vg view, & vg gampcompare (Comparing graph alignments for mapping benchmark) 096bfdce

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 and reviewers. 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

All data used in this study are available at <https://github.com/jonassibbesen/vgrna-project-paper> (DOI: 10.5281/zenodo.7234532). It links the following public data resources:

GRCh38 reference genome (primary assembly) [ftp://ftp.ensembl.org/pub/release-94/fasta/homo\\_sapiens/dna/Homo\\_sapiens.GRCh38.dna.primary\\_assembly.fa.gz](ftp://ftp.ensembl.org/pub/release-94/fasta/homo_sapiens/dna/Homo_sapiens.GRCh38.dna.primary_assembly.fa.gz)  
 GENCODE v29 (primary assembly) [ftp://ftp.ebi.ac.uk/pub/databases/genocode/Genocode\\_human/release\\_29/genocode.v29.primary\\_assembly.annotation.gtf.gz](ftp://ftp.ebi.ac.uk/pub/databases/genocode/Genocode_human/release_29/genocode.v29.primary_assembly.annotation.gtf.gz)  
 1000 Genomes Project phased variant calls [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/GRCh38\\_positions/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/GRCh38_positions/)  
 IPD-IMGT/HLA (release 3.43.0) <https://github.com/ANHIG/IMGTHLA>

And also the simulated data used for this study:

Simulated reads: [http://cgl.gi.ucsc.edu/data/vgrna/simulated\\_data/](http://cgl.gi.ucsc.edu/data/vgrna/simulated_data/)

Downsampled GENCODE annotation [http://cgl.gi.ucsc.edu/data/vgrna/transcript\\_annotation/](http://cgl.gi.ucsc.edu/data/vgrna/transcript_annotation/)

And sequencing reads (by accession or by sample identifier, depending on the data repository):

ENCSR000AED replicate 1 (ENCODE)

ENCSR706ANY (ENCODE)

ENCSR146ZKR (ENCODE)

ENCSR825GWD (ENCODE)

ENCSR686JJB (ENCODE)

ENCSR502OTI (ENCODE)

ENCSR995BHD (ENCODE)

SRR1153470 (SRA)

ERR1050073 (SRA)

ERR1050074 (SRA)

ERR1050075 (SRA)

ERR1050076 (SRA)

ERR1050077 (SRA)

ERR1050078 (SRA)

ERR1050079 (SRA)

ERR1050080 (SRA)

ERR1050081 (SRA)

SRR12765650 (SRA)

SRR12765534 (SRA)

CHM13 replicate 1 (T2T)

NA07051, NA11832, NA11840, NA11930, NA12287, NA12775, NA12889, NA06994, NA07037, NA07357, NA11829, NA11893, NA12006, NA12043, NA12234, NA12272 & NA12275 (Geuvadis)

<https://github.com/jonassibbesen/vgrna-project-paper> also includes links to spliced pangenome graphs and pantranscriptomes constructed in this study, which may be freely used in other projects.

Mapping benchmark tables and haplotype-specific expression estimates are available at DOI: 10.5281/zenodo.7234454

## Field-specific reporting

Please select the one below that is 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/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

# Life sciences study design

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

|                 |   |
|-----------------|---|
| Sample size     | All real sequencing data sets used in this study were obtained from public repositories, so the sample size was determined by the coverage available from these sources. The sample sizes of all simulated sequencing data sets were uniformly chosen to be 50,000,000. This sample size is similar to the sequencing coverage available from the public data sets. Additionally, it is sufficiently large to estimate mapping accuracy for a mapper with an error rate of $10^{-6}$ , which corresponds to the conventional maximum mapping quality of Phred 60. |
| Data exclusions | No data was excluded from analyses. Results for data that were used to optimize the method are included as supplementary figures.   |
| Replication     | The results shown in the main figures are from data sets that were never used in the process of optimizing the methods. All optimizations were carried out on separate data sets, and the analyses were then repeated on a held out data set. The results on the training data sets are included as supplementary figures. Broadly speaking, the results on the held out data sets and training data sets are similar.  |
| Randomization   | This study does not use any randomized study designs. The objects of comparison are different computational methods, not different treatments or conditions. The computational methods are all applied to the exact same data, so there is no need to use randomization to control for unobserved correlates.   |
| Blinding        | Blinding is not relevant to this study. There were no human participants, and all analyses were fully determined and automated before analyzing the held-out data sets.   |

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

| n/a                                 | Involved in the study                                  |
|-------------------------------------|--|
| <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 and archaeology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Human research participants   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Dual use research of concern  |

### 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 |