



Long-read reconstruction of many diverse haplotypes with devider

Jim Shaw, Christina Boucher, Yun William Yu, et al.

Genome Res. published online September 23, 2025
Access the most recent version at doi:[10.1101/gr.280510.125](https://doi.org/10.1101/gr.280510.125)

P<P	Published online September 23, 2025 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Long-read reconstruction of many diverse haplotypes with devider

Jim Shaw^{*1,2}, Christina Boucher³, Yun William Yu⁴, Noelle Noyes⁵, and Heng Li^{1,2}

¹ Department of Data Science, Dana-Farber Cancer Institute, Boston, MA 02215, USA

² Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02215, USA

³ Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL 32611, USA

⁴ Ray and Stephanie Lane Computational Biology Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA

⁵ Department of Veterinary Population Medicine, University of Minnesota, St. Paul, MN 55421, USA

*Corresponding author: jshaw@ds.dfci.harvard.edu

Abstract. Reconstructing exact haplotypes is important when sequencing a mixture of similar sequences. Long-read sequencing can connect distant alleles to disentangle similar haplotypes, but handling sequencing errors requires specialized techniques. We present devider, an algorithm for haplotyping small sequences—such as viruses or genes—from long-read sequencing. devider uses a positional de Bruijn graph with sequence-to-graph alignment on an alphabet of informative alleles to provide a fast assembly-inspired approach compatible with various long-read sequencing technologies. On a synthetic Oxford Nanopore Technologies (ONT) long-read dataset containing seven HIV strains, devider recovered 97% of the haplotype content and had the most accurate abundance estimates while taking < 4 minutes and 1 GB of memory for > 8000× coverage. Benchmarking on synthetic mixtures of antimicrobial resistance (AMR) genes showed that devider recovered 83% of haplotypes, 23 percentage points higher than the next best method. On real PacBio and ONT datasets, devider recapitulates previously known results in seconds, disentangling a bacterial community with > 10 strains and an HIV-1 co-infection dataset. We used devider to investigate the within-host diversity of a long-read bovine gut metagenome enriched for AMR genes, discovering 13 distinct haplotypes for a *tet(Q)* tetracycline resistance gene with > 18,000× coverage and 6 haplotypes for a *CfxA2* beta-lactamase gene. We found clear recombination blocks for these AMR gene haplotypes, showcasing devider’s ability to unveil evolutionary signals for heterogeneous mixtures.

Keywords: long-reads · haplotyping · viruses · genes · metagenome · de Bruijn graph

1 Introduction

2 The presence of highly similar genomic sequences within a single or a group of organisms is common in
3 biological settings. Examples include viral quasispecies (Domingo and Perales, 2019) in single-stranded RNA
4 virus populations (e.g., HIV-1 and SARS-CoV-2) (Cuevas et al., 2015) or co-existing microbial subspecies
5 in microbiomes (Van Rossum et al., 2020). Small genomic differences can have large functional implica-
6 tions (Vedantam et al., 1998; Olkkola et al., 2010), so it is crucial to disentangle this heterogeneity. We
7 will call the recovery process of similar genomic sequences “haplotyping” or “phase”. Although traditionally
8 used in the context of diploid organisms, we extend the concept here to encompass the resolution of genetic
9 diversity in microbes, viruses, or even genes.

10 With high-throughput sequencing, we can obtain haplotypes by linking reads that share informative al-
11 leles, for example, single nucleotide polymorphisms (SNPs). Unfortunately, standard de novo short-read or
12 long-read assembly approaches can collapse small-scale variation (Bickhart et al., 2022), returning only a
13 consensus sequence. Although haplotype-resolved assembly has become standard for PacBio HiFi sequenc-
14 ing (Cheng et al., 2021; Feng et al., 2022; Benoit et al., 2024; Li and Durbin, 2024), HiFi data may not always
15 be available and assembly is computationally intensive. In contrast to assembly approaches, reference-based
16 haplotyping uses a reference plus alignment to facilitate haplotyping; many existing approaches use the align-
17 ment, SNP calling, then phasing paradigm (Feng et al., 2021; Knyazev et al., 2021; Cai and Sun, 2022; Cai
18 et al., 2022; Edge et al., 2017; Shaw and Yu, 2022; Lancia et al., 2001; Zhou et al., 2024; Patterson et al.,
19 2015).

20 We are interested in reference-based haplotyping for (1) long-read sequencing, (2) small sequences of
21 approximately the read length, and (3) an *unknown*, possibly large number of haplotypes. Reference-free
22 approaches that tackle all or a subset of criteria (1-3) also exist (Luo et al., 2022; Baaijens et al., 2019), but
23 the lack of reference adds additional algorithmic difficulties; we focus on the reference-based case. Long reads
24 can connect more distant alleles across shared genomic regions than short reads. Still, a technical challenge is
25 to deal with sequencing errors for certain technologies, e.g., Oxford Nanopore Technologies (ONT) long reads
26 can have 90 – 99% sequencing accuracy depending on the chemistry and basecalling (Sereika et al., 2022).
27 We focus on small sequences on the order of the read length (i.e., “local haplotyping” (Zagordi et al., 2011)),
28 but we do not necessarily require all reads to overlap the region of interest. This is sufficient for haplotyping
29 genes of interest or estimating diversity. Although this seems like a simple task, systematic errors, high
30 coverage, and low abundance haplotypes make accurate reconstruction challenging (Eliseev et al., 2020).
31 Another class of approaches focuses on genome-scale haplotype reconstruction for prokaryotes (Shaw et al.,

2024; Vicedomini et al., 2021; Kazantseva et al., 2024). However, these methods have not been tested for phasing high-diversity communities with > 5 similar haplotypes.

Many haplotyping algorithms work for long reads or an unknown number of haplotypes, but only a few were designed to do both. RVHaplo (Cai and Sun, 2022), HaploDMF (Cai et al., 2022), CliqueSNV (Knyazev et al., 2021), and iGDA (Feng et al., 2021) are long-read methods that call low-frequency SNPs and can phase diverse genomic sequences; RVHaplo uses a network clustering formulation, HaploDMF uses a matrix factorization approach, CliqueSNV uses a clique-merging approach, and iGDA uses a probabilistic local haplotyping step with an overlap-layout algorithm for reconstruction.

We present *devider*, a new long-read, reference-based haplotyping method for diverse small sequences. Given a set of aligned reads and SNPs, *devider* models the haplotyping problem as an assembly problem on a positional de Bruijn graph (PDBG). *devider* is inspired by the kSNP algorithm (Zhou et al., 2024) which uses a PDBG similarly but only for haplotyping diploids. We use the fact that the PDBG naturally splits if enough variation is present and collapses under ambiguity, thus haplotyping samples without prior knowledge of the number of distinct sequences. We then leverage the length of long reads by finding walks along the PDBG through unitig construction and read-to-graph alignment. We find that *devider* efficiently resolves haplotypes in a variety of synthetic and real datasets and show its versatility to reveal genomic heterogeneity.

Methods

At a high level, *devider* takes a set of aligned reads in BAM format plus SNPs in VCF format. It then outputs the sequences of SNPs for each recovered sequence (i.e., the haplotype), the abundance of each haplotype, and a base-level sequence for each haplotype. *devider* does not call SNPs or perform alignment, but we supply a wrapper for *devider* with *minimap2* (Li, 2018) and *LoFreq* (Wilm et al., 2012) for alignment and SNP calling. We use *LoFreq* for variant calling unless otherwise stated, but users can choose to use their own SNP caller.

devider follows in three major steps (**Fig. 1A-C**): (1) encoding the aligned reads with SNPs, (2) constructing a PDBG and positional unitig graph while filtering errors, and (3) aligning the SNP-encoded reads to the positional unitig graph to obtain walks along the graph that are supported by many reads, which represent candidate haplotypes.

60 **Encoding reads with informative SNPs and filtering false positives**

61 The first step is to filter false positive SNPs arising from strand-specific errors. For each SNP, we apply
62 a simple strand bias filter using Fisher’s exact test (Guo et al., 2012) on the 2-by-2 contingency table
63 defined by forward/reverse strands and reference/alternate alleles. For high coverage, Fisher’s exact test
64 is too stringent, so we also require an odds ratio of > 1.5 or $< 1/1.5$ for the contingency table. We then
65 apply the Benjamini-Hochberg (Benjamini and Hochberg, 1995) multiple-testing correction at 0.005 FDR
66 threshold. More powerful methods are possible for strand filtering (McElroy et al., 2013), but we opt for a
67 simple method because we are only concerned with phase.

68 We then subsample the SNPs if too many are present (i.e. the sequences are highly divergent) as follows:
69 for a given parameter α (discussed in **Practical details and implementation**) and median number of
70 SNPs contained in a read β , we downsample uniformly to retain only α/β of the SNPs if $\alpha < \beta$. We do this
71 because when too many SNPs are present, a group of SNPs may span only a small region of the genome,
72 which we wish to avoid in our subsequent SNP-based k -mer approach. We then realign each read using
73 blockaligner (Liu and Steinegger, 2021) against the 32 bp flanks around each filtered SNP site, replacing the
74 site with all possible alleles and then selecting the allele that gives the highest alignment score.

75 Finally, we encode each read as an ordered list of tuples such as (3, 1), (4, 1), (5, 0), (6, 1), (7, 0), (8,0);
76 see **Fig. 1A**. The first number represents the i -th SNP in the reference to which the read is aligned, and the
77 second number indicates the allele that the read contains where 0 is the reference and 1 to 3 indicate the
78 alternates. We skip SNPs if either the SNP site has a deletion in the read or the read’s base at the SNP site
79 is neither an alternate nor reference allele in the VCF. Thus, a sequence such as (3, 1), (5, 0) is possible.

80 **Positional de Bruijn graph and unitigs on the SNP alphabet**

81 All reads will now be considered SNP-encoded reads. An SNP encoded read is a string in an alphabet
82 $\Sigma = \mathbb{Z}^+ \times \{0, 1, 2, 3\}$ subject to the constraint that the first symbol in each letter is in increasing order, e.g.,
83 (5,1) must come before (7,0). k -mers of SNP-encoded reads are defined as elements in Σ^k , and we construct
84 a de Bruijn graph in the usual way by collecting all k -mers within the reads and adding directed edges
85 between k -mers that overlap $k - 1$ SNPs (**Fig. 1B**, top). However, because the positions are encoded in
86 the alphabet Σ , we only collapse k -mers if they have the same positions. This is now a *positional* de Bruijn
87 graph (PDBG) (Ronen et al., 2012; Bao et al., 2014; Cameron et al., 2017). The PDBG is a directed acyclic
88 graph (DAG) since any cycle would violate the increasing ordering of the SNPs. This fact will be important
89 for the subsequent sequence-to-graph alignment steps.

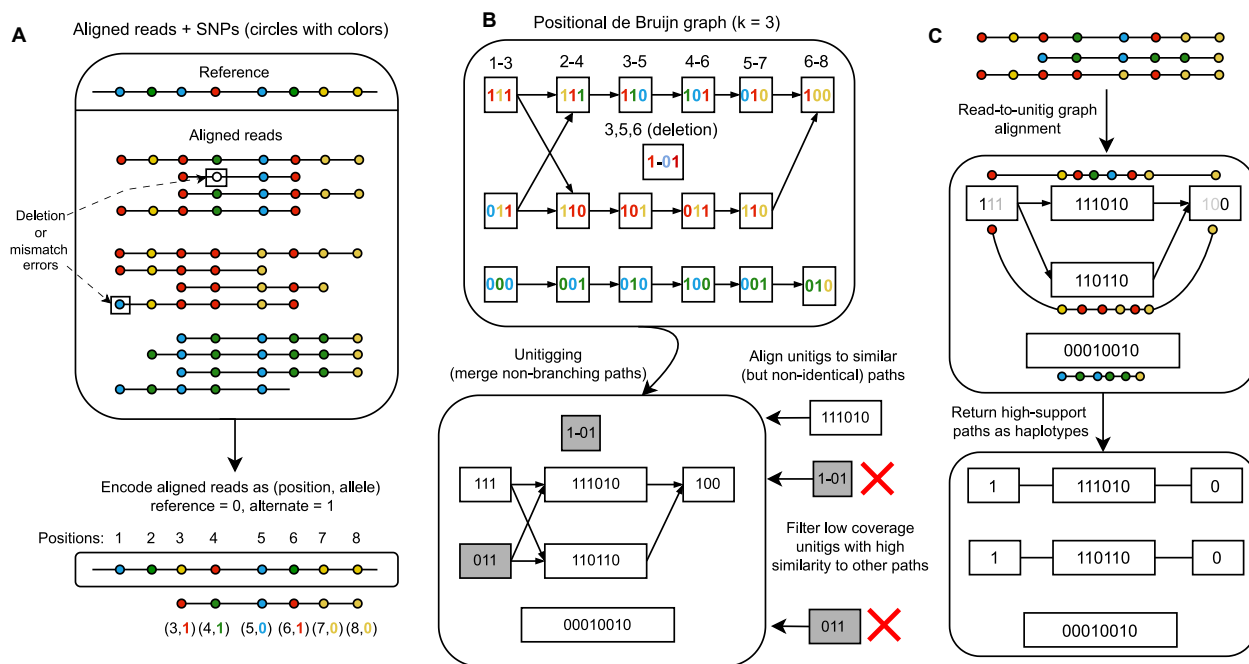


Fig. 1. Algorithmic framework for devider. **A.** Reads that are aligned to a reference are converted to a SNP representation with positional information. Sequencing errors lead to erroneous SNP encodings. **B.** The SNP-encoded reads are turned into a positional de Bruijn graph (PDBG) ($k = 3$ shown). In a PDBG, k -mers are collapsed if their alleles *and* their positions are identical. Errors in reads lead to spurious k -mers in the PDBG. After merging paths with in-degree and out-degree equal to 1 (unitigging), unitigs are aligned back to the graph to filter low-coverage, high-similarity unitigs. **C.** Reads are aligned back to the filtered unitig graph to determine high-confidence walks through the graph. These paths are taken to be putative haplotypes. devider then post-processes the haplotypes to output haplotype abundances, a base-level consensus of each haplotype, and the reads assigned to each haplotype.

90 We automatically choose a value of k as follows. Let γ be the 33rd percentile of the number of SNPs
 91 contained in a read, and let N be the number of SNPs in the reference (after filtering). Let M be a parameter
 92 representing the maximum possible value of k , which is set to avoid long error-prone k -mers. This depends
 93 on the sequencing technology and is picked through a preset option. We let $k = \min(M, N \cdot \frac{3}{4}, \gamma)$. We do not
 94 let k span more than 75% of the reference to not miss k -mers if a smaller haplotype only covers a subsection
 95 of the reference. We discuss the parameter choices of M in **Practical details and implementations** and
 96 show results over varying values of k in the **Results**.

97 We apply an initial filtering step to discard the likely erroneous k -mers. Let the coverage of a k -mer be
 98 the number of times it appears in a SNP-encoded read, and let m be the mean k -mer coverage. Let A be the
 99 minimum allowable abundance for a haplotype (default = 0.0025 or 0.25%). We filter k -mers that appear
 100 only once or have less than $m \cdot A$ coverage. Finally, from the filtered PDBG, we construct the positional
 101 unitig graph by merging all non-branching paths into unitigs (**Fig. 1B**, bottom) (Myers, 2005). We let the
 102 coverage of the unitig be the mean coverage of the merged k -mers.

103 **Filtering unitigs by error-aware unitig-to-graph alignment**

104 The main technical challenge is to simplify the unitig graph, which can still have many spurious unitigs
 105 arising from k -mer errors (**Supplementary Figure 1 and 2**). In the standard de Bruijn graph assembly,
 106 tip removal and bubble popping are used to remove noise *and* variation (Li et al., 2015). However, we do
 107 not want to remove the true variation. Thus, we use a unitig-to-graph alignment approach plus coverage
 108 information for fine-grained unitig filtering; this generalizes both tip removal and bubble popping but uses
 109 alignment information and coverage information.

110 **Classifying errors** Let G be the positional unitig graph. Recall that a unitig node v can be represented by
 111 $v = (x_1, x_2, \dots, x_n)$ where $x_i = (a_i, b_i)$ with a_i the SNP position and b_i the allele. Given two unitigs v_1 and v_2 ,
 112 we classify errors as SNP deletions (*del*), reference-to-alternate (*rtoa*) mismatches, and alternate-to-reference
 113 (*ator*) mismatches as follows.

- 114 – $s(v_1, v_2)$ is the number of SNPs that are the same in v_1 and v_2 (i.e., share the same position and base).
- 115 – $del(v_1, v_2)$ is the number of SNPs in v_1 are deleted relative to v_2 (i.e., do not appear in v_2 and lie between
 116 the first and last SNP of v_2).
- 117 – $rtoa(v_1, v_2)$ and $ator(v_1, v_2)$ represent the number of SNPs that have the same position but different
 118 alleles between v_1 and v_2 . Specifically, $rtoa(v_1, v_2)$ is the number of different SNPs for which v_1 has the
 119 *reference* allele (and thus v_2 has an alternate allele), while $ator$ is the number of differing SNPs for which
 120 v_1 has an *alternate* allele.

121 We stratify the error types because they appear with different frequencies. *rtoa* and *ator* differ due to
 122 reference bias: If a read comes from a haplotype with a true alternate allele, it may systematically align
 123 with the reference allele incorrectly (Stevenson et al., 2013). A SNP deletion (*del*) error is the most common
 124 due to the following reason: Consider a biallelic site with two alleles (A, C) and a read originating from
 125 the haplotype with allele A. divider’s convention is to consider the read’s SNP to be deleted if there is a
 126 base-level deletion in the CIGAR string or the read’s base is G or T at the SNP site. Of the four error
 127 possibilities (deletion, substitution to C, substitution to G, substitution to T), three of them result in a SNP
 128 deletion in the SNP-encoded reads. Furthermore, long reads can also have inherently higher deletion error
 129 rates (Delahaye and Nicolas, 2021).

130 **Alignment to DAG** Next, we align the unitigs back to the unitig graph to find an alignment path. However,
 131 we disallow alignments back to unitig itself, since this would always be the best match. If there is a path

132 such that s (successes) is large relative to $del, rtoa$ and $ator$ (errors) and this path has much higher coverage
 133 than the unitig, then the unitig is probably an error that originates from the path.

134 By abuse of notation, let $s(v, P)$ refer to the number of matching alleles between v and the string spelled
 135 out by a unitig path P . We wish to find a path that does not contain v and maximizes $s(v, P)$. Furthermore,
 136 we impose that (1) ties are broken by taking the highest coverage path, (2) v 's first and last SNP positions
 137 must be within P 's first and last SNP positions, and (3) all unitigs in P overlap v . We do not need to penalize
 138 for the error terms in this step because more matching alleles imply fewer deletions and errors.

Due to the DAG structure of the PDBG (and thus the positional *unitig* graph), the optimal path can be found with standard dynamic programming. Taking a unitig v and a topological order on nodes $1, \dots, n$ that overlap v , we wish to find the optimal path ending at node v_j . Let P_j be a path that ends at node v_j . The following recurrence holds:

$$\text{score}(j) := \max_{P_j} s(v, P_j) = \max_{v_\ell \in \text{in}(v_j)} \left[\max_{P_\ell} s(v, P_\ell) + s(v, v_j) - s(v, \text{overlap}(v_j, v_\ell)) \right].$$

139 Note that we subtract the overlap to avoid double counting, since v_j and its incoming unitigs may overlap.
 140 After obtaining all scores, we find optimal paths that satisfy the three constraints above. Each alignment
 141 takes $O(|V||E|)$ worst-case running time, but in practice, unitig graphs are sparse (e.g. **Supplementary**
 142 **Fig. 1**), so this step is not a bottleneck.

143 **Error-aware filtering** Once we have a best path P for each unitig v , we use a one-sided binomial test
 144 based on coverage to filter spurious unitigs. This filtering is “error-aware” in the sense that different types
 145 of errors in the sequence-to-graph alignment have different frequencies (see section **Classifying errors**), so
 146 we wish to differentiate the types of errors within the binomial test. We run the alignment and error-aware
 147 filter starting from the smallest coverage unitig, then remove unitigs from the graph, and repeat for the next
 148 smallest coverage unitig.

149 Let $cov(v)$ be the unitig coverage, and let $cov(P)$ be (1) the highest coverage unitig in P that completely
 150 covers v or (2) the mean unitig coverage in P if no unitigs in P completely cover v . We filter out v if
 151 $\Pr(cov(v) \geq \text{Binomial}(cov(P), q)) < 0.005$, where q is defined as follows: define $p_{del} = 0.35$, $p_{rtoa} = 0.15$,
 152 and $p_{ator} = 0.10$ to model the error frequency as previously discussed. Then $q = p_{ator}^{ator(v,P)} \cdot p_{rtoa}^{rtoa(v,P)} \cdot \delta(p_{del})$
 153 where $\delta(p_{del}) = p_{del}$ if $ator(v, P) + rtoa(v, P) = 0$ but otherwise equal to 1. This formula slightly deviates
 154 from the independence assumption of the three error modalities. We found that systematic biases could
 155 occasionally cause two errors to occur non-independently. Thus, we loosen the independence of $p(del)$, which
 156 was the main cause of errors. In general, we found that these values are conservative and work for error-prone

reads ($\approx 95\%$), but could likely be tightened to improve sensitivity in the future as sequencing technologies improve.

Read-to-graph alignment and collecting haplotype paths

To find candidate haplotype paths through the filtered unitigs, we align reads back to the graph using a sequence-to-DAG alignment algorithm (**Fig. 1C**), which is similar to unitig-to-graph alignment. We will use these alignments to find well-supported walks through the PDBG, representing candidate haplotypes.

Let r be a read. Here, rather than finding a path to maximize $s(r, P)$ as we did before, we find a path to maximize $s(r, P) - 3 \times err(r, P)$ where $err = ator(r, P) + rtoa(r, P)$. The penalty term of 3 was chosen to penalize errors more than similarities. Although true deleted SNPs in a haplotype are possible, we assume deletions occur mainly due to noise, so we do not penalize deletions. Compared to unitig-to-graph alignment, we add two changes. First, we allow P to bridge sinks-to-sources. This allows the alignment to rescue broken paths due to erroneous filtering. Secondly, we do not require the path to completely cover the read, also in case of erroneous filtering. We use the same dynamic programming procedure as for unitig-to-graph alignment; $err(r, P)$ can also be split into the exact same recurrence and is thus solvable by dynamic programming.

After aligning all reads, we have a set of paths with read-alignment multiplicities. If two or more equally good paths exist for a read, we do not assign the read to any path. If a path is contained within another, we remove the contained path. If the path is *uniquely* contained in another, we add the contained path's multiplicity to the non-contained path's multiplicity. We filter for high-confidence paths by removing paths for which the read-alignment multiplicity is $< 3 \times$ the minimum unitig coverage within the path.

Haplotype consensus and outputs For the filtered paths, we assign the reads to each path by finding the path that maximizes $2 * s(r, P) - 3 * rtoa(r, P) - 5 * ator(r, P) - del(r, P)$, assigning the read to no path if the score is less than 0. The penalty weights are heuristically chosen to reflect the frequency of occurrence: deletions are common so they are penalized less, and $rtoa(r, P)$ is penalized less than $ator(r, P)$ because of reference bias; the read r will preferentially carry the reference allele due to biases in aligning to references. We then take the consensus allele for each SNP site after assignment and the abundance as the fraction of assigned reads.

Lastly, we perform a deduplication step as follows. For some resolution parameter ρ , we merge two resulting consensus haplotypes H_1, H_2 together if $(rtoa^\Delta(H_1, H_2) + ator^\Delta(H_1, H_2)) / s^\Delta(H_1, H_2) < \rho$, where Δ considers only unambiguous SNPs for which the > 0.75 fraction of reads carry the majority allele. After merging, we reassign the reads to the best candidate haplotype subject to the same scoring scheme above, take consensus, and filter low-abundance and low-depth haplotypes (default = 0.25% and $5 \times$ depth). We

188 iterate this procedure until the number of haplotypes does not change. We then return the abundances,
189 the reads assigned to each haplotype, and the sequences of the SNPs for each haplotype. Finally, we also
190 output a majority *base-level* consensus sequence of all bases, not just SNPs, for each set of assigned reads
191 by iterating through alignments in the BAM file. We return the ‘N’ base if the fraction of reads supporting
192 the majority base is less than a parameter (default = 0.66).

193 **Practical details and implementation**

194 `devider` is implemented in Rust and uses `rust-htslib` (Bonfield et al., 2021) and `rust-bio crates` (Köster, 2016).
195 We implemented the following preset options to help users pick parameters: `old-long-reads`, `nanopore-r9`,
196 `nanopore-r10`, `hi-fi`. These correspond to $(M, \alpha, \rho) = (10, 50, 0.02)$, $(20, 150, 0.01)$, $(35, 250, 0.005)$, and
197 $(100, 500, 0.001)$ respectively. The current default is `nanopore-r9`. We wrap `devider` in `minimap2` and `LoFreq`
198 in the `run_devider_pipeline` script in the repository. This runs `minimap2`, `SAMtools`, and `LoFreq` to gener-
199 ate an indexed BAM and VCF pair. For `LoFreq`, we found that disabling base-alignment qualities with the
200 `-B` option improved sensitivity on nanopore reads. All other parameters were set to default.

201 **Results**

202 We show our benchmarking setups in **Supplementary Figure 3 and 4**. For a set of genomes, we use
203 `badread` (Wick, 2019) to simulate nanopore long reads with default settings except for lengths and accu-
204 racy, which we make explicit for each data set in the following. We picked an arbitrary reference genome
205 and ran all tools with this reference genome. We then compared their predicted haplotypes with the true
206 haplotypes. We benchmarked `devider` against `iGDA v1.0.1`, `CliqueSNV v2.0.3`, `RVHaplo v2`, and `Hap-`
207 `loDMF` (version May 2022). We tried running `Strainline`, a de novo viral quasi-species assembly method,
208 but it produced an error that was related to an unresolved GitHub issue ([https://github.com/HaploKit/](https://github.com/HaploKit/Strainline/issues/17)
209 `Strainline/issues/17`). We ran all methods with default settings and `iGDA` in its ONT setting with the
210 `ont_context_effect_read_qv_12_base_qv_12` model. We ran all the methods with 10 threads. We used
211 the following metrics for benchmarking:

- 212 1. Hamming SNP error: the mean percentage of incorrect SNPs for each predicted haplotype against its
213 best-matching genome.
- 214 2. Fraction recovered: the fraction of SNPs recovered for all genomes after matching predicted haplotypes
215 against their best-matching genomes.
- 216 3. Earth mover’s distance (EMD) (Rubner et al., 1998): a measure of the distance between the predicted
217 abundances and the true abundances. The pairwise distance function we used for the EMD is the number

218 of mismatched SNPs between the predicted and true haplotype, and the weights are the predicted relative
 219 abundances.

220 4. Haplotype error: the predicted number of haplotypes minus the true number.

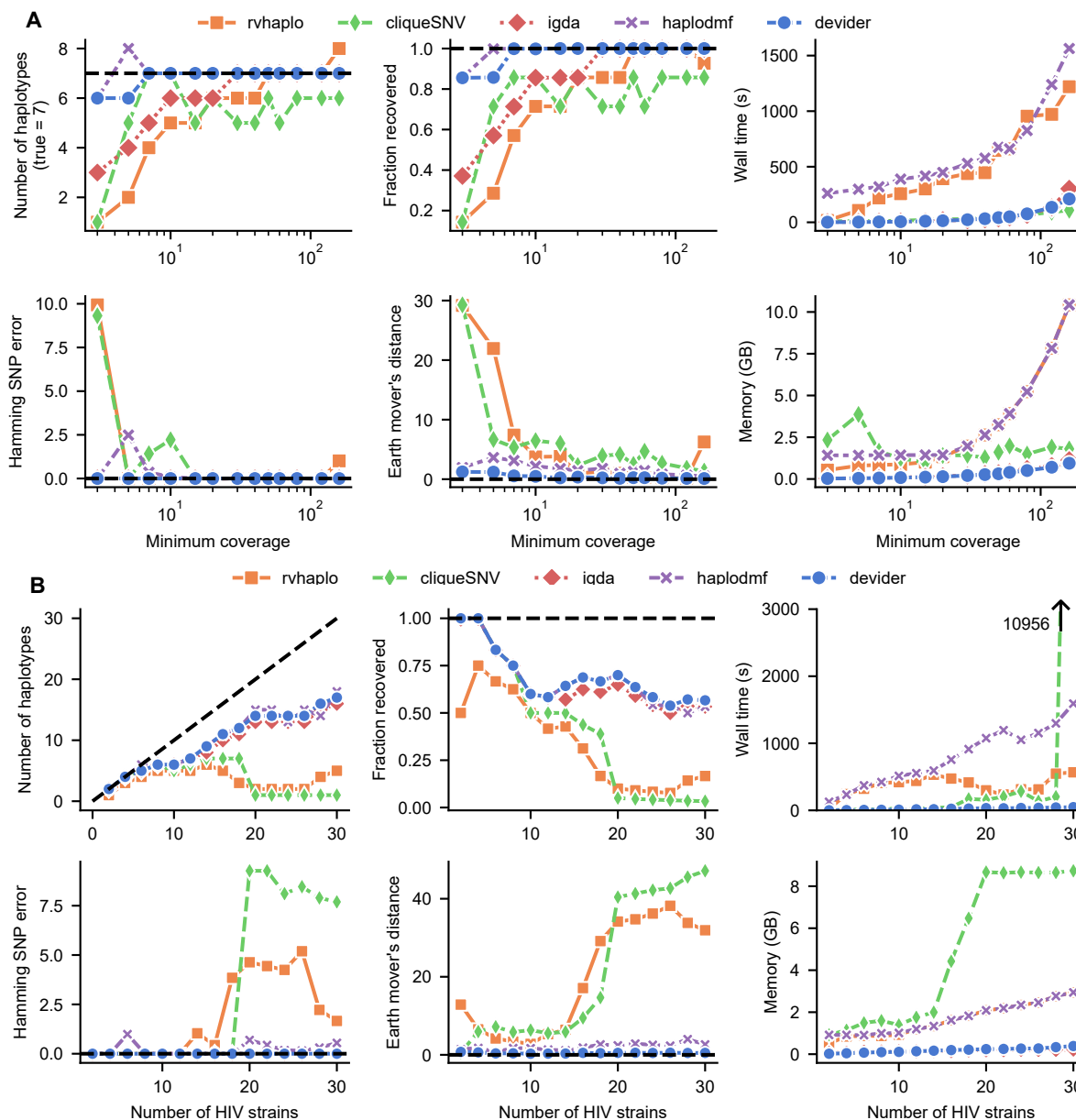


Fig. 2. Benchmarking long-read haplotyping tools on simulated HIV-1 communities. **A.** 7 HIV-1 strains from Kinloch et al. (Kinloch et al., 2023) at staggered abundances (1:3:5:7:9:10:20) with simulated reads (9000 bp mean length; 95% mean accuracy). The x-axis indicates the depth of coverage for the lowest-coverage strain. **B.** 2 to 30 HIV-1 strains with $15\times$ - $145\times$ uniformly random coverage and the same simulated read lengths/accuracy. CliqueSNV failed to complete within the default 3-hour time limit for the last sample. *igDA does not output abundances, so the earth mover's distance could not be computed. Dashed black lines indicate optimal performance

221 HIV-1 benchmarking (7 strains plus varying coverage and 2-30 strains)

222 HIV-1 serves as a standard benchmarking genome for viral quasispecies methods due to its fast mutation rate
223 as an ssRNA virus and its clinical and public health importance. Thus, we created two HIV-1 communities:
224 a 7-strain staggered abundance community and multiple communities from 2 to 30 strains with uniform
225 abundances. These strains ranged from 99.27% to 99.71% pairwise nucleotide similarity. Here, we define
226 “strain” to be a distinct reference genome that we are trying to reconstruct.

227 **7 strains at staggered abundances** We took a set of 30 HIV-1 genomes from Kinloch et al. (Kinloch
228 et al., 2023) with accessions available in Supplementary Table 1. For the first 7 strain dataset, we selected
229 an arbitrary reference (OR483991.1) and the 7 most similar strains as determined by skani (Shaw and Yu,
230 2023). The abundances were staggered at a 1:3:5:7:9:10:20 ratio, with the smallest strain coverage ranging
231 from $3\times$ to $160\times$. We simulated reads in three settings with (Accuracy, Length) = (95%, 9000bp), (98%,
232 9000bp), and (95%, 3000bp) with length standard deviation of 500bp. The precision of 95% represents
233 older or faster nanopore sequencing runs, while 98% is more representative of the best current basecall-
234 ing/chemistries (Sereika et al., 2022). HIV-1 genomes are approximately 9000 bp, so the two length settings
235 represent complete and partial coverage.

236 On the 95% accuracy, 9000bp dataset (**Fig. 2A**), *devider* and HaploDMF performed the best. *devider* had
237 slightly worse mean fraction recovered (97.7% for *devider* vs 98.8% for HaploDMF), equal mean haplotype
238 error (-0.15 for *devider* vs +0.15 for HaploDMF), and better EMD (0.41 for *devider* vs 1.66 for HaploDMF).
239 However, *devider* was > 13 times faster and took < 10 times less memory than HaploDMF on average. Only
240 *devider* and *iGDA* achieved perfect SNP Hamming errors across all data points. The greatest performance
241 difference was at low coverage; at $3\times$ minimum coverage, *devider* estimated 6 correct haplotypes and Hap-
242 loDMF estimated 8 (incorrectly outputting an additional haplotype), but the other methods only recovered
243 1, 1, and 3 haplotypes for RVHaplo, CliqueSNV, and *iGDA* respectively. We found that CliqueSNV consis-
244 tently obtained 6 haplotypes, missing the low abundance haplotype. We tried to increase its sensitivity by
245 lowering its abundance threshold to 0.25% (the same as *devider*), but it drastically overestimated the number
246 of haplotypes, outputting > 30 haplotypes at high coverage. To show that *devider* is robust to parameter
247 choices on this dataset, we varied the k -mer length between 10 to 30 and found that *devider* still had the
248 best haplotype error, fraction recovered, and EMD (**Supplementary Fig. 5**).

249 On the 95%, 3000bp dataset (**Supplementary Fig. 6 a**), *devider* had the second best mean fraction
250 recovered (92.3% vs 93.7% for HaploDMF) and the best mean EMD (1.18 vs 4.18 for the second-best
251 HaploDMF). *devider* had the smallest absolute mean haplotype error (0.46 vs 0.53 for the second-best

252 HaploDMF) compared to the other methods. However, the mean Hamming SNP error for devider (0.37%)
253 was slightly worse than iGDA (0.06%) and RVHaplo (0.07%), but better than HaploDMF (1.02%).

254 On the 98% accuracy, 9000bp dataset (**Supplementary Fig. 6 b**), devider and HaploDMF had near-
255 perfect performances for mean fraction recovered (98.9% vs 99.9% respectively), haplotype error (-0.08 vs
256 -0.15 respectively). devider had perfect Hamming SNP error (0% vs 0.024% for iGDA, the second best) and
257 the best EMD (0.31 vs 2.3 for HaploDMF, the second best). de Bruijn graph methods work well with lower
258 error rates because the probability of having an error within a k -mer is approximately $(accuracy)^k$; this is
259 encouraging for future nanopore datasets as read accuracy increases.

260 **2-30 strains at uniform coverages** We simulated reads at 95% accuracy and 9000 bp average length
261 for 2-30 strains from the same set of HIV-1 genomes, with each genome at $15\times$ - $145\times$ coverage uniformly at
262 random (**Fig. 2B**). devider, iGDA, and HaploDMF performed better than RVHaplo and CliqueSNV on this
263 dataset. devider was the best method on all metrics except HaploDMF recovered slightly more haplotypes
264 on average (10.2 vs 10.1 for devider). While no method could capture all strains when > 10 were present,
265 devider consistently had low EMD, implying that missed strains were either lowly abundant or highly similar.
266 Unlike HaploDMF, RVHaplo, and CliqueSNV, devider had a lower EMD as the number of strains increased.
267 In fact, devider had a smaller EMD with 30 strains than with two strains. CliqueSNV also performed well
268 when the number of strains was < 10 (81.7% mean fraction recovered compared to 83.6% for devider), but
269 its performance dropped when more than 10 strains were present (37.5% mean fraction recovered for 10-20
270 strains versus 65.6% for devider).

271 devider and iGDA stood out in terms of efficiency compared to the other methods. iGDA and devider
272 took 5 and 23 seconds on average, respectively, and < 0.5 GB of RAM. Note that we included lofreq's runtime
273 and memory usage in devider's results. RVHaplo took 370 seconds, HaploDMF took 790, and CliqueSNV
274 820 seconds on average. CliqueSNV was efficient except for the case with 30 strains, where the runtime
275 ballooned to 10,956 seconds, which corresponds to approximately 3.04 hours. We found that this was because
276 CliqueSNV sets a 3 hour (10,800 seconds) time limit by default if it cannot solve the haplotyping problem,
277 after which it outputs no haplotypes.

278

279 **SARS-CoV-2 minor haplotyping**

280 We next investigated the ability of algorithms to detect minor haplotypes at uneven abundances. To do this,
281 we created multiple two-strain synthetic mixtures of Delta (accession MZ009823.1) and Omicron (accession
282 OL672836.1) SARS-CoV-2 genomes, with the minor Omicron strain ranging from 0.39% to 25% abundance.

283 This setup represents a situation where we would like to detect a circulating, low-abundance strain with high
284 sequence similarity (99.58% between these two genomes) prior to emergence. This was exactly the case with
285 SARS-CoV-2 evolution after the Delta phase of the pandemic: a pattern of selective sweeps and emergence
286 of a single low-abundance strain characterized viral strain frequencies (Markov et al., 2023; Boyle et al.,
287 2022).

288 We first simulated near full-length reads (30kb mean length with 500bp standard deviation) and 95%
289 accuracy at $3000\times$ total coverage, i.e., a 1% abundant haplotype would have $30\times$ coverage (**Fig. 3A**).
290 devider and RVHaplo were the best two methods in this dataset, successfully constructing exactly two haplotypes in
291 $3/8$ and $5/8$ of the mixtures respectively. The limit of detection of RVHaplo was 1.56%, slightly better than
292 devider at 3.12% and devider incorrectly output a spurious low abundance haplotype at 25% abundance.
293 iGDA did not estimate two haplotypes in any abundances for this dataset, HaploDMF output incorrect
294 numbers of haplotypes except at 25% abundance, and CliqueSNV failed to detect the minor haplotype below
295 12.5% abundance.

296 In practice, capturing the entire SARS-CoV-2 in a single read is difficult, and approaches focus on
297 amplicon sequencing of shorter target regions (Tyson et al., 2020). For a more realistic test, we focus on the
298 spike protein region of SARS-CoV-2, a smaller region of approximately 4kbp for which there are long-read
299 amplicon protocols (Nimsamer et al., 2023; Liao et al., 2022). We simulated reads for the spike protein
300 region in the same genomes at $3000\times$ total coverage but now with 4kb mean length and 500bp standard
301 deviation (**Fig. 3B**). In this setup, devider and RVHaplo were again the best two methods, but this time
302 having almost identical performance at a detection limit of 3.12%. HaploDMF was able to recover the minor
303 haplotype at 1.56% abundance, but HaploDMF also outputs an additional spurious haplotype. In general,
304 detecting haplotypes of very low abundance in abundances $< 2\%$ with 95% accuracy reads is a difficult task,
305 and RVHaplo was the best method with devider in second.

306 **Synthetic AMR gene haplotyping for 53 AMR gene groups**

307 We chose AMR genes as a gene-level haplotyping benchmark due to their diversity in sequence composition,
308 length, biological significance, and also the prevalence of targeted enrichment sequencing protocols (Slizovskiy
309 et al., 2022; Shay et al., 2023; Baba et al., 2023) for which our methods are potentially usable. We clustered
310 AMR genes in the MEGARes database version 3.0 (Bonin et al., 2023) by computing all pairwise average
311 nucleotide identities (ANI) and mean alignment fractions (AF) with skani (using the `--slow` preset) and
312 then using the Leiden algorithm (Traag et al., 2019; Camargo et al., 2024) with edge weights as $ANI * AF$
313 at resolution 1.00. Of the remaining 53 clusters with ≥ 15 different haplotypes, we sampled between 2-15
314 haplotypes at coverage $80\times$ - $1000\times$, both uniformly at random. We then simulated reads at 95% accuracy

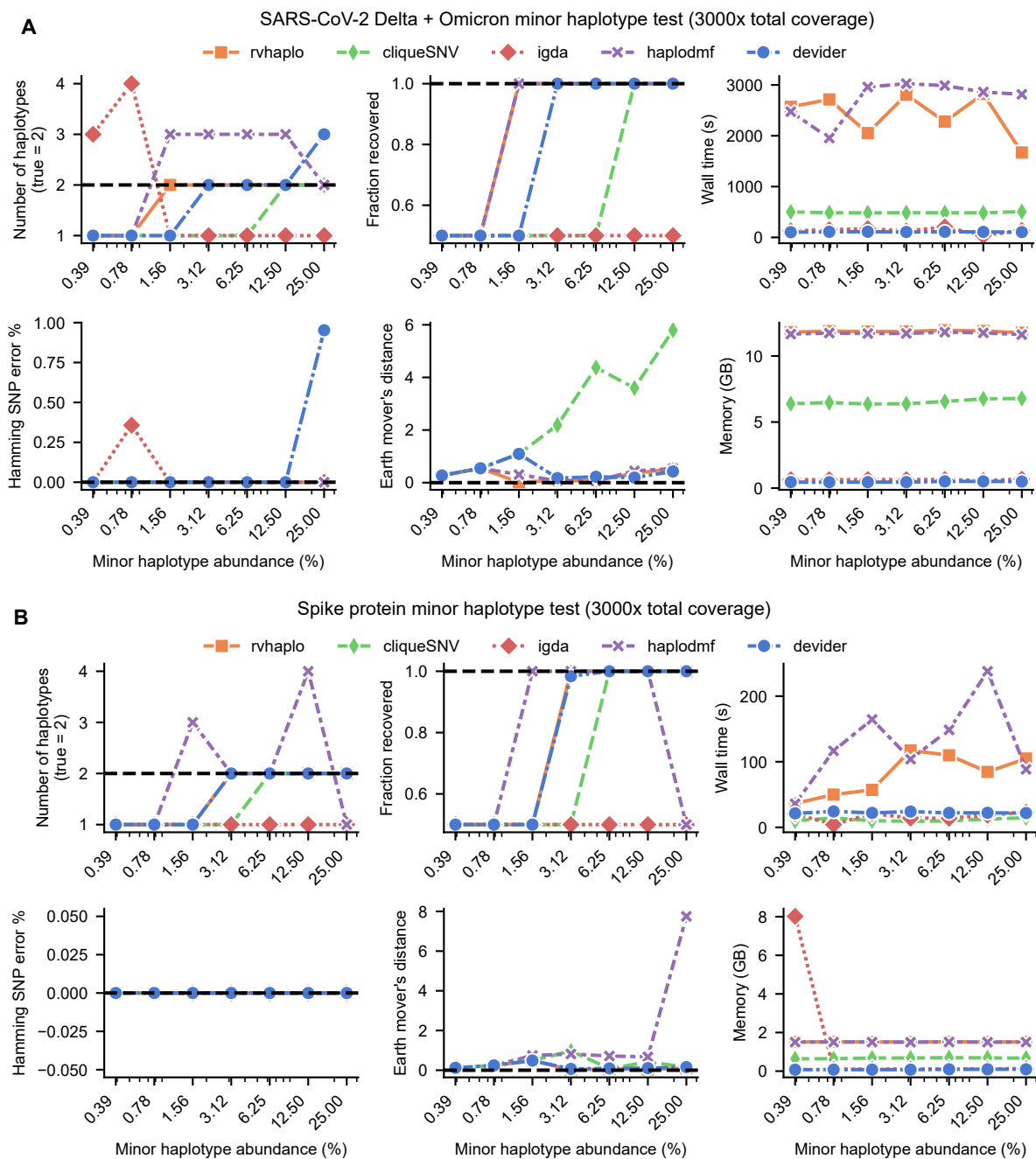


Fig. 3. Benchmarking for two-strain SARS-CoV-2 synthetic sequencing mixtures at varying abundances. **A.** Results for mixtures of full-length reads of Delta (major) and Omicron (minor) genomes. **B.** Results for mixtures of reads that only cover the spike protein gene for the Delta and Omicron genomes.

315 and 1500 bp length with 200bp standard deviation. The lengths of the AMR genes ranged from 721 to 3303
 316 bp.

317 In this data set, devider was the best method for the mean haplotype error, SNP Hamming error, and
 318 fraction recovered (**Fig. 4A**), with -1.5 , 0.06% and 83.3% , respectively. CliqueSNV was the next best

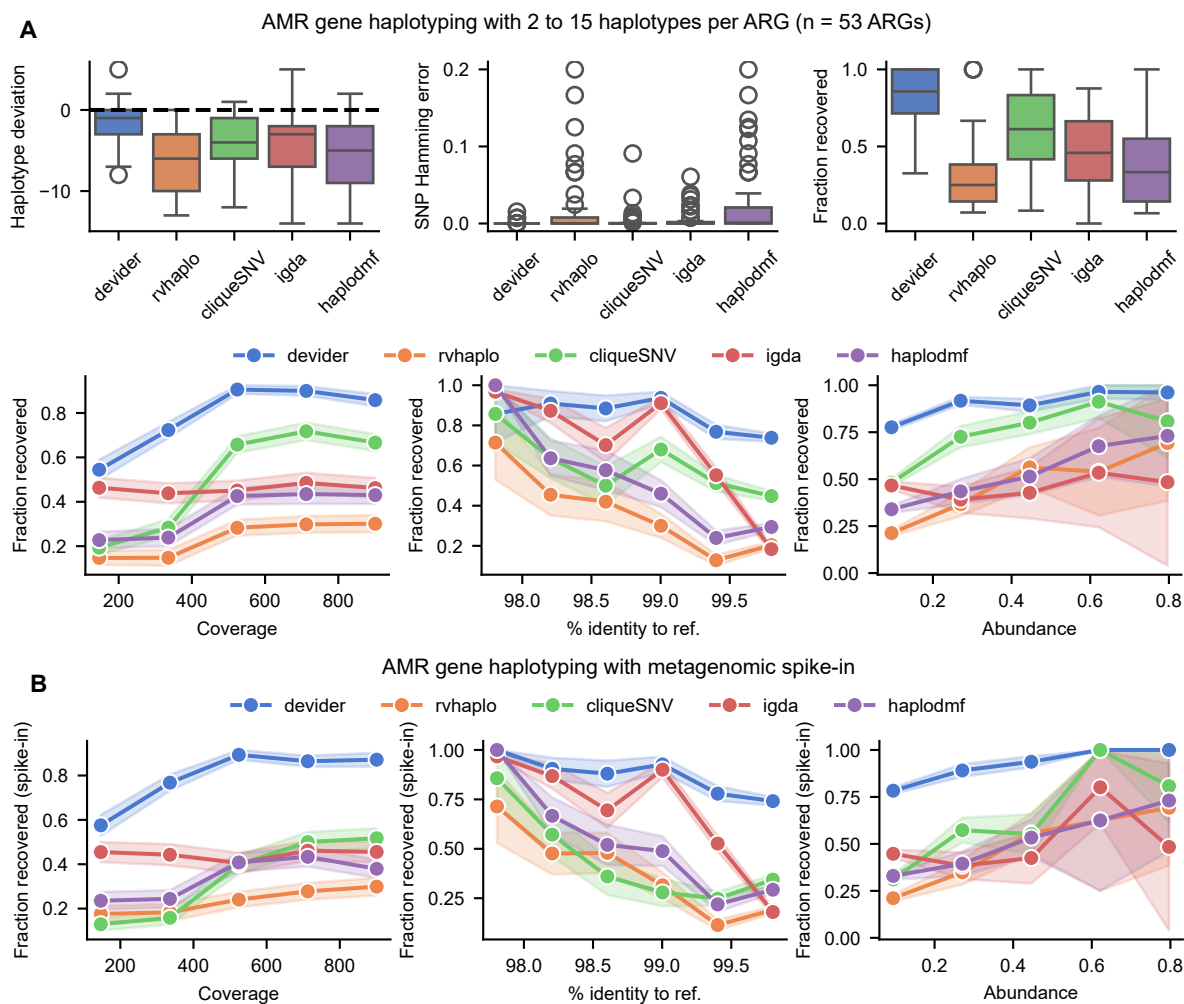


Fig. 4. Benchmarking for haplotyping of synthetic mixtures of antimicrobial resistance (AMR) genes. **A.** Top: results for 53 sets of AMR genes with 2-15 haplotypes and reads simulated at $80\times - 1000\times$ coverage (95% identity; 1500 mean length), both picked uniformly at random. Bottom: the same results but with fraction recovered as a function of the haplotype's coverage, its % nucleotide identity to the reference, and its abundance (i.e., normalized coverage). **B.** Fraction recovered for the AMR genes after spiking the AMR reads into a synthetic long-read mouse gut metagenome from CAMI2. Each method was rerun after aligning the pooled dataset against the AMR gene references. Error bars indicate standard errors after binning data points along the x-axis. Box plots show the median, the 25th and 75th percentiles, and $1.5\times$ the interquartile range.

319 method with -4.2 , 0.34% , and 60.2% on the same metrics. To investigate the discrepancy between different
 320 methods, we stratified the fraction recovered by coverage, % identity of haplotype to reference, and abundance
 321 (**Fig. 4B**). As expected, lower coverage, which is correlated with lower abundance, leads to lower fraction
 322 recovered for all methods. We found that the likely cause of the performance discrepancy was high similarity
 323 haplotypes: $> 2/3$ of the haplotypes had $> 99.5\%$ similarity to the reference, and in these haplotypes, the
 324 recovered fraction of devider was 73.8% compared to CliqueSNV with 46.3% , the next best method. Thus,

devider is the most sensitive method for highly similar haplotypes, which is important because capturing even small variations in AMR genes can alter phenotypes (Vedantam et al., 1998).

AMR dataset with spike-in metagenome We extended the previous AMR haplotyping experiment to a metagenomics setting. In metagenomics, the input is a mixture of microbial genomic reads *and* reads from AMR genes. Thus, we mixed the previously simulated AMR reads into a simulated long-read mouse gut metagenome from CAMI2 (Meyer et al., 2022) (labeled as sample 0), hereafter referred to as the spiked metagenome.

In particular, genomes within the CAMI2 metagenome contain AMR genes or sequences with homology to AMR genes. Mapping the CAMI2 reads against MEGARES with minimap2 resulted in 130 AMR genes with $> 2\times$ coverage but only 10 AMR genes with $> 20\times$ coverage. Thus, the spike-in metagenome contains the previously simulated AMR haplotypes (with 2-15 haplotypes and $80\text{-}1000\times$ coverage) as well as this new tail of low-abundance AMR haplotypes from the CAMI2 metagenome.

In this setup (shown in **Supplementary Fig. 4**), we measured the ability of each method to recover the original simulated AMR haplotypes (**Fig. 4B**). devider, iGDA, HaploDMF, and RVHaplo had a $< 2\%$ difference in fraction recovered compared to the previous (without spike-in) case. However, CliqueSNV fell to 42.9% (with spike-in) from 60.2% (no spike-in). Thus, devider can recover abundant haplotypes for a long-tailed abundance distribution with low abundance haplotypes, a common characteristic of metagenomics data.

Results on real heterogeneous sequencing mixtures

HIV-1 co-infection haplotyping Mori et al. (Mori et al., 2022) sequenced a set of HIV-1 samples with full-length nanopore amplicons (5.8-7% sequencing error rate) and detected a possible HIV-1 co-infection for a patient (labeled TRN9) by stochastically subsampling reads and reassembling. Here, we ran devider, iGDA, RVHaplo, and CliqueSNV on these reads and an arbitrary reference (NC_001802.1) to try to confirm their results. devider gave four haplotypes (60.0%, 32.34%, 4.12%, and 3.57% abundance) and RVHaplo gave two (67% and 32.9% abundance). CliqueSNV only output one haplotype and iGDA outputs no haplotypes.

Mori et al. found two majority haplotypes, which is corroborated by RVHaplo and devider. We show devider's two majority haplotypes in **Fig. 5A**. We investigated the two additional minority haplotypes output by devider to see if they were erroneous; curiously, these two haplotypes were a mix of the two majority haplotypes and well supported by "chimeric" reads with noisy breakpoints (**Supplementary Fig. 7**). HIV-1 is known to recombine, which may be a possible explanation, but in vitro PCR recombination could also be a possibility (Meyerhans et al., 1990). Porechop (Wick et al., 2017) found only 6/2500 reads

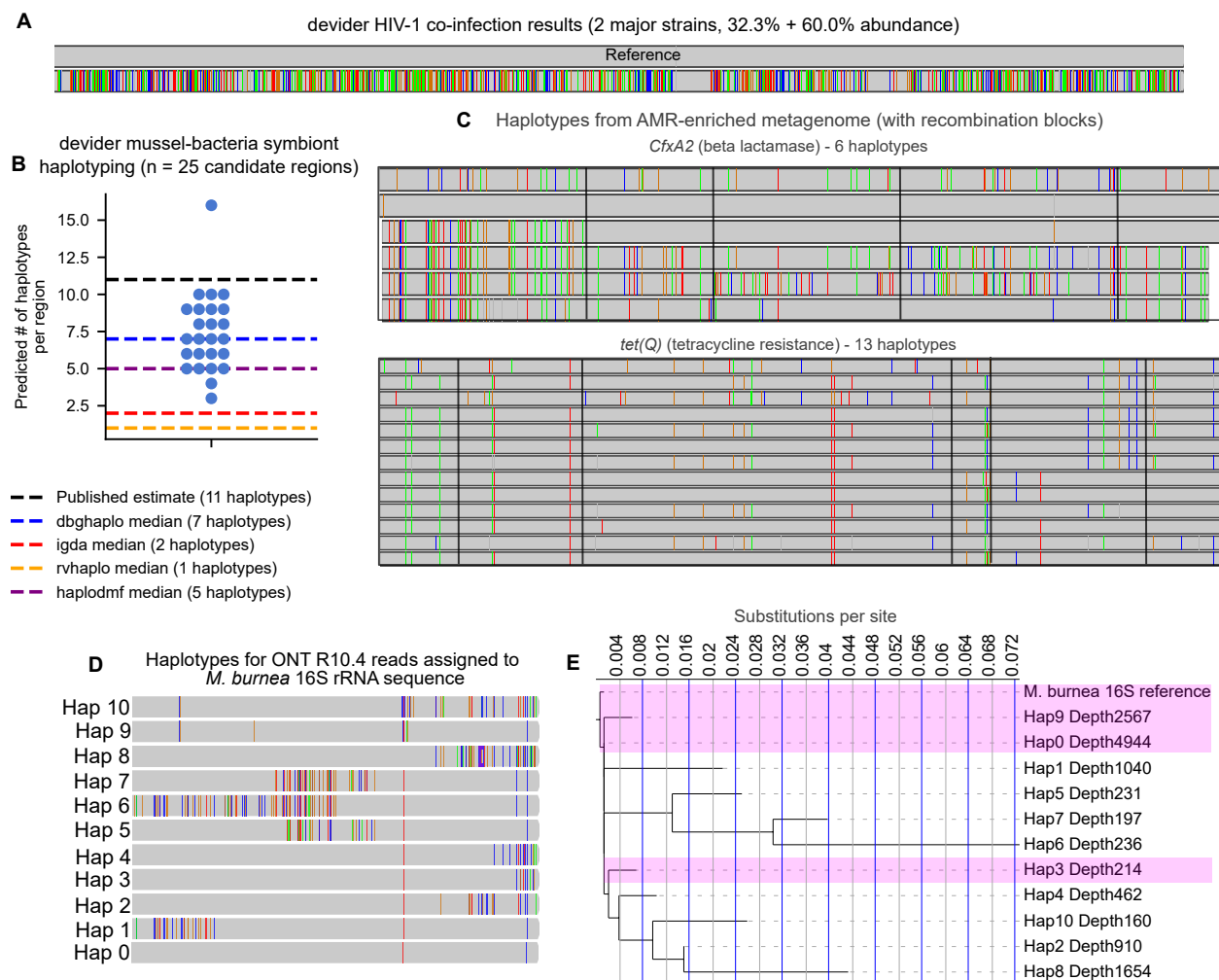


Fig. 5. Long-read haplotyping results from real samples subjected to a variety of sequencing technologies. **A.** Haplotypes from long-read HIV-1 nanopore sequencing (93 - 94.2% predicted sequencing accuracy) of an HIV co-infection from Mori et al. (Mori et al., 2022). Two major haplotypes were found by devider, confirming previous results. Mismatched bases are shown with the reference as the upper haplotype. **B.** Haplotyping results for PacBio RS II sequencing (89.5% mean gap-compressed identity against reference) of an intracellular bacterial symbiont community within deep-sea mussels from Ansorge et al. (Ansorge et al., 2019), who predicted 11 strains to be present. 25 candidate single-copy regions with high SNP diversity were haplotyped by devider, iGDA, and RVHaplo; CliqueSNV was excluded because it timed out on multiple regions. devider produced higher diversity estimates than iGDA and RVHaplo, which both produced ≤ 3 haplotypes across all sites. **C.** devider haplotyping of a long-read bovine gut metagenome enriched for AMR genes. *CfxA2* (3200 \times coverage) and *tet(Q)* (19500 \times coverage and last 1000bp shown) haplotype sequences with $> 30\times$ coverage and 1% abundance are shown with mismatches against their reference sequences in MEGARes v3.0. Mismatches shared by all haplotypes are removed. Recombination blocks are outlined in black as predicted by GARD. **D.** devider haplotypes from an ONT R10.4 16S rRNA dataset for the reference 16S sequence of the most abundant species *M. burnea*. **E.** Phylogenetic tree of haplotypes assigned to the *M. burnea* reference. Depth of coverage is shown next to the haplotype ID. The x-axis shows the branch length from the root. Highlighted haplotypes have $> 99\%$ identity to the reference.

356 with an adapter in the middle, so an erroneous ligation is not likely. We do not speculate further on the
 357 biology, but we note that these chimeric reads represent a real signal that devider can retrieve for further
 358 biological investigation.

359 **Mussel-bacteria symbiont community with 11 estimated strains** Ansorge et al. (Ansorge et al.,
360 2019) investigated the strain-level diversity of an intracellular, sulfur-oxidizing bacterial symbiont community
361 within deep-sea mussels. They sequenced one sample using the PacBio RS II, which produces high-error
362 long reads. Ansorge et al. do not explicitly define “strain”, but calculate it via two methods: (1) counting
363 unique structural arrangements of single-copy marker genes from long-read assemblies and (2) using a viral
364 quasispecies detection algorithm (Jayasundara et al., 2015) for short-read sequencing in the same dataset. We
365 investigated whether long-read haplotyping without assembly could give similar estimates, even with noisy
366 reads (mean gap-compressed identity against reference = 89.5%). To generate reasonable diversity estimates,
367 we haplotyped 3kb regions of the reference genome (GCF_900128535) that had > 10 SNPs (as detected by
368 LoFreq) and coverage between 200 and 250. These criteria were chosen because the average read length was
369 3.8 kbp and the estimated single-copy coverage, calculated by dividing the total read bases by the genome
370 size, was 220. We ran all the methods on these regions in PacBio mode (if such a preset existed) and the
371 `old-long-reads` preset for devider.

372 In these 25 regions, devider estimated a median of seven haplotypes (**Fig. 5B**). RVHaplo estimated a
373 median of five haplotypes, but iGDA and RVHaplo only estimated ≤ 2 median haplotypes. CliqueSNV failed
374 to run on multiple regions, timing out after three hours and outputting nothing, so we did not include
375 its results. devider estimated 16 haplotypes for one region, and subsequent investigation revealed many
376 additional alignments near the edge of contigs. The first 600 bases had mean coverage > 300 \times , possibly
377 indicating a duplicated region (**Supplementary Fig. 8**) for a subset of the strains. Ultimately, devider was
378 able to capture some of the known diversity in these samples using noisy long reads, whereas most other
379 methods failed.

380 **Discovering recombinant AMR genes in AMR-enriched long-read metagenomes** We used devider
381 to haplotype an AMR-enriched PacBio CCS long-read fecal metagenome from a cow that received intensive
382 antibiotic treatments (Slizovskiy et al., 2022). We first dereplicated the MEGARes v3.0 database at 95%
383 using vsearch (Rognes et al., 2016) to avoid ambiguous mapping of reads to highly similar genes. We then
384 ran devider with extra stringent parameters, setting the minimum abundance to 1% and minimum depth to
385 30 \times .

386 In total, we found 18 different dereplicated AMR genes with ≥ 2 haplotypes. Of these 18 genes, 9 were
387 tetracycline resistance genes that were phased into 52 distinct haplotypes. The highest coverage gene was
388 *tet(Q)* (Lacroix and Walker, 1996) at 19500 \times coverage and phased into 13 haplotypes. Other high-diversity
389 genes we found included *mefA* (Daly et al., 2004), a macrolide efflux pump, at 7100 \times coverage and phased
390 into 12 haplotypes, as well as *CfxA2* (Iwahara et al., 2006), a beta-lactamase, at 3200 \times coverage and phased

391 into 6 haplotypes. We illustrate the haplotypes of *tet(Q)* and *CfxA2* in **Fig. 5C** (IGV (Robinson et al.,
392 2011) screenshots in **Supplementary Fig. 9 and 10**). We found a distinct mosaic structure within these
393 haplotypes, suggesting a history of recombination within these haplotypes. We used MAFFT (Kato et al.,
394 2002) to generate a multiple sequence alignment from devider's haplotypes and GARD (Kosakovsky Pond
395 et al., 2006) to detect recombination, which found evidence of recombination for both genes. We draw
396 breakpoints where GARD's model-averaged support was > 0.3 in **Fig. 5C**. The consensus haplotypes were
397 well-supported by the reads: the 6 haplotypes for *CfxA2* had 68%, 72%, 90%, 71%, 44%, and 70% of their
398 assigned reads spanning all 4 recombination breakpoints (in order from top to bottom of **Fig. 5C**). Across
399 all alleles, a median of 99% of the reads within the haplotype supported the consensus allele, indicating
400 confident haplotypes. Mosaicism due to recombination is a well-documented characteristic of some ribosomal
401 protection proteins including *tet(Q)* (Warburton et al., 2016), and *CfxA* genes are commonly colocalized
402 with an element known to play a role in the mosaic behavior of conjugative elements (García et al., 2008).
403 Thus, these detected recombination events are supported by known mechanisms in these two AMR genes.

404
Disentangling 16S rRNA amplicon sequences from R10.4 nanopore data As an additional use case,
405 we investigated using devider as a reference-based method to cluster full-length 16S rRNA amplicon sequences
406 from ONT R10.4 sequencing, the newest and most accurate chemistry. Currently, computational profilers for
407 ONT 16S sequencing align amplicons directly to reference genomes (Curry et al., 2022). Denoising algorithms
408 for generating amplicon sequencing variants (ASVs) require a significant fraction of error-free reads (Callahan
409 et al., 2016), so they are still not usable for the newest R10.4 reads. We investigated whether devider clusters
410 can be used to generate reference-based ASVs, allowing for species-level identification. We profiled a 16S
411 R10.4 soil sample from Zhang et al. (2023) (accession SRR23176498) as follows. We first used Emu (Curry
412 et al., 2022) to quantify species-level abundances and then applied devider with Emu's default 16S database.
413 We used the same pipeline for devider except parameters `-B 2 -A 3 -s 20` for minimap2 (more aggressive
414 extension) and `-mapq-cutoff 1 -supp-mapq-cutoff 1 -min-qual 20` for devider (MAPQ is low for 16S
415 databases; higher base quality thresholds for newer nanopore reads). devider only took 200 seconds for the
416 entire dataset, excluding alignment and variant calling.

417
418 For the most abundant species, *Massilia eburnea* (18.6% abundance from Emu), we visualized the 11
419 haplotypes found by devider and a phylogenetic tree (Letunic and Bork, 2021) of the consensus sequences
420 constructed by FastTree (Price et al., 2010) and MAFFT (**Fig. 5D, E**). Each consensus had $> 90x$ depth of
421 coverage and were well-supported by reads (**Supplementary Fig. 11**). Only Hap0, Hap3, and Hap9 had $>$
422 99% identity to the reference, a suggested threshold for species-level assignment (Edgar, 2018). These three

423 haplotypes had 59.9% combined abundance, so using this 99% identity threshold, 40.1% of the sequences
424 should be considered novel species-level 16S sequences. As a more extreme example, *Vicinamibacter silvestris*
425 was the third most abundant species according to Emu (4.5% abundance); however, every one of the devider
426 consensus had < 95% identity to the reference (**Supplementary Fig. 12**). Ultimately, robustly generating
427 ASVs is a highly non-trivial task that we do not claim to solve. However, our investigation shows how
428 devider can be a useful tool for curating reference-based ASVs from deep, heterogeneous long-read amplicon
429 sequences.

430 Discussion

431 We presented devider, a method for retrieving high-similarity haplotypes from long-read sequencing of het-
432 erogeneous sequences. devider leverages a positional de Bruijn graph (PDBG) assembly approach on a subset
433 of informative alleles to disentangle variation. This framework is efficient and naturally resolves variation
434 without the need to explicitly infer the number of haplotypes. The key technical challenge was to remove
435 sequencing artifacts within the PDBG, especially for error-prone long reads, while retaining high sensitivity,
436 which we accomplished through an error-aware sequence-to-graph alignment approach.

437 Based on our benchmarks and experience, we found devider to excel for heterogeneous and high depth
438 samples. Key examples include amplicon sequencing, enriched metagenome sequencing, or viral sequenc-
439 ing. Anecdotally, we have found devider to also work for high-abundance species in unenriched long-read
440 metagenomes. In general, devider can work for < 10× depth (**Fig. 2A**), but its sensitivity increases for
441 higher depth. Currently, we do not try to recover haplotypes with less than 0.25% abundance and < 5×
442 depth by default. In practice, the exact detection limits will be some function of relative abundance, depth,
443 and sequence divergence (see **Fig. 4**). We have shown that devider can distinguish up to ≈ 20 distinct hap-
444 lotypes in benchmarks and real data, although more could be possible depending on the relative divergence
445 of haplotypes.

446 We designed devider to work with a wider range of technologies and sequencing error rates. We limited
447 devider to reconstructing “small” sequences on the order of read length for conservative recovery. However,
448 we showed that synthetic reconstructing an HIV genome of even > 3 times the mean read length is possible
449 as long as some of the reads are long enough. As error rates improve, it may be possible to attempt a longer
450 haplotype reconstruction using our approach.

451 A key limitation is our reference-based approach, which is unable to recover new sequences de novo.
452 However, reference-based approaches are intrinsically more efficient and simpler than de novo approaches.
453 We believe that reference-based methods are complementary to de novo approaches. As the capabilities and

20 J. Shaw et al.

454 the need for resolved sequences at the haplotype level from long reads continue to increase, *devider* will be
455 a fast and useful tool for retrieving accurate haplotypes.

456 **Software availability**

457 *devider* is open source and is available on GitHub (<https://github.com/bluenote-1577/devider>) or bio-
458 conda (Grüning et al., 2018) and as Supplemental Code. The scripts for reproducing our figures are available
459 at <https://github.com/bluenote-1577/devider-test> and as Supplemental Scripts.

460 **Competing interest statement**

461 The authors declare no conflict of interest.

462 **Acknowledgments**

463 N.N. is supported by the National Institutes of Health grant 1R01AI173928-01A1. C.B. is supported by
464 the National Institutes of Health grant 5R01AI141810-02. Y.W.Y. is supported by the US National Science
465 Foundation grant 2531433. H.L. is supported by the National Institutes of Health grant R01HG010040. J.S.
466 is supported by a Natural Sciences and Engineering Research Council of Canada Postdoctoral Fellowship.

467 Author contributions: C.B., Y.W.Y, N.N, and J.S. conceived the study. J.S. designed and implemented
468 the method with supervision from H.L. J.S. performed the bioinformatics analysis and benchmarking with
469 supervision from all authors. All authors contributed to the writing of the manuscript.

Bibliography

- 471 Ansorge R, Romano S, Sayavedra L, Porras MÁG, Kupczok A, Tegetmeyer HE, Dubilier N, and Petersen
472 J. 2019. Functional diversity enables multiple symbiont strains to coexist in deep-sea mussels. *Nature*
473 *Microbiology* **4**: 2487–2497.
- 474 Baaijens JA, Van der Roest B, Köster J, Stougie L, and Schönhuth A. 2019. Full-length de novo viral
475 quasispecies assembly through variation graph construction. *Bioinformatics* **35**: 5086–5094.
- 476 Baba H, Kuroda M, Sekizuka T, and Kanamori H. 2023. Highly sensitive detection of antimicrobial resistance
477 genes in hospital wastewater using the multiplex hybrid capture target enrichment. *mSphere* **8**: e00100–23.
- 478 Bao E, Jiang T, and Girke T. 2014. AlignGraph: Algorithm for secondary de novo genome assembly guided
479 by closely related references. *Bioinformatics* **30**: i319–i328.
- 480 Benjamini Y and Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful Approach
481 to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)* **57**: 289–300.
- 482 Benoit G, Raguideau S, James R, Phillippy AM, Chikhi R, and Quince C. 2024. High-quality metagenome
483 assembly from long accurate reads with metaMDBG. *Nature Biotechnology* .
- 484 Bickhart DM, Kolmogorov M, Tseng E, Portik DM, Korobeynikov A, Tolstoganov I, Uritskiy G, Liachko I,
485 Sullivan ST, Shin SB, et al.. 2022. Generating lineage-resolved, complete metagenome-assembled genomes
486 from complex microbial communities. *Nature Biotechnology* **40**: 711–719.
- 487 Bonfield JK, Marshall J, Danecek P, Li H, Ohan V, Whitwham A, Keane T, and Davies RM. 2021. HTSlib:
488 C library for reading/writing high-throughput sequencing data. *GigaScience* **10**: giab007.
- 489 Bonin N, Doster E, Worley H, Pinnell LJ, Bravo JE, Ferm P, Marini S, Prospero M, Noyes N, Morley PS,
490 et al.. 2023. MEGARes and AMR++, v3.0: An updated comprehensive database of antimicrobial resis-
491 tance determinants and an improved software pipeline for classification using high-throughput sequencing.
492 *Nucleic Acids Research* **51**: D744–D752.
- 493 Boyle L, Hletko S, Huang J, Lee J, Pallod G, Tung HR, and Durrett R. 2022. Selective sweeps in SARS-CoV-2
494 variant competition. *Proceedings of the National Academy of Sciences* **119**: e2213879119.
- 495 Cai D, Shang J, and Sun Y. 2022. HaploDMF: Viral haplotype reconstruction from long reads via deep
496 matrix factorization. *Bioinformatics* **38**: 5360–5367.
- 497 Cai D and Sun Y. 2022. Reconstructing viral haplotypes using long reads. *Bioinformatics* **38**: 2127–2134.
- 498 Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, and Holmes SP. 2016. DADA2: High-
499 resolution sample inference from Illumina amplicon data. *Nature Methods* **13**: 581–583.

- 500 Camargo AP, Call L, Roux S, Nayfach S, Huntemann M, Palaniappan K, Ratner A, Chu K, Mukherjee S,
501 Reddy TBK, et al.. 2024. IMG/PR: A database of plasmids from genomes and metagenomes with rich
502 annotations and metadata. *Nucleic Acids Research* **52**: D164–D173.
- 503 Cameron DL, Schröder J, Penington JS, Do H, Molania R, Dobrovic A, Speed TP, and Papenfuss AT.
504 2017. GRIDSS: Sensitive and specific genomic rearrangement detection using positional de Bruijn graph
505 assembly. *Genome Research* **27**: 2050.
- 506 Cheng H, Concepcion GT, Feng X, Zhang H, and Li H. 2021. Haplotype-resolved de novo assembly using
507 phased assembly graphs with hifiasm. *Nature Methods* **18**: 170–175.
- 508 Cuevas JM, Geller R, Garijo R, López-Aldeguer J, and Sanjuán R. 2015. Extremely High Mutation Rate of
509 HIV-1 In Vivo. *PLoS Biology* **13**: e1002251.
- 510 Curry KD, Wang Q, Nute MG, Tyshaieva A, Reeves E, Soriano S, Wu Q, Graeber E, Finzer P, Mendling W,
511 et al.. 2022. Emu: Species-level microbial community profiling of full-length 16S rRNA Oxford Nanopore
512 sequencing data. *Nature Methods* **19**: 845–853.
- 513 Daly MM, Doktor S, Flamm R, and Shortridge D. 2004. Characterization and Prevalence of MefA, MefE,
514 and the Associated msr(D) Gene in *Streptococcus pneumoniae* Clinical Isolates. *Journal of Clinical*
515 *Microbiology* **42**: 3570–3574.
- 516 Delahaye C and Nicolas J. 2021. Sequencing DNA with nanopores: Troubles and biases. *PLOS ONE* **16**:
517 e0257521.
- 518 Domingo E and Perales C. 2019. Viral quasispecies. *PLoS Genetics* **15**: e1008271.
- 519 Edgar RC. 2018. Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *Bioinformatics* **34**:
520 2371–2375.
- 521 Edge P, Bafna V, and Bansal V. 2017. HapCUT2: Robust and accurate haplotype assembly for diverse
522 sequencing technologies. *Genome Research* **27**: 801–812.
- 523 Eliseev A, Gibson KM, Avdeyev P, Novik D, Bendall ML, Pérez-Losada M, Alexeev N, and Crandall KA.
524 2020. Evaluation of haplotype callers for next-generation sequencing of viruses. *Infection, genetics and*
525 *evolution : journal of molecular epidemiology and evolutionary genetics in infectious diseases* **82**: 104277.
- 526 Feng X, Cheng H, Portik D, and Li H. 2022. Metagenome assembly of high-fidelity long reads with hifiasm-
527 meta. *Nature Methods* **19**: 671–674.
- 528 Feng Z, Clemente JC, Wong B, and Schadt EE. 2021. Detecting and phasing minor single-nucleotide variants
529 from long-read sequencing data. *Nature Communications* **12**: 3032.
- 530 García N, Gutiérrez G, Lorenzo M, García JE, Píriz S, and Quesada A. 2008. Genetic determinants for cfxA
531 expression in *Bacteroides* strains isolated from human infections. *Journal of Antimicrobial Chemotherapy*
532 **62**: 942–947.

- 533 Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, Valieris R, and Köster J. 2018.
534 Bioconda: Sustainable and comprehensive software distribution for the life sciences. *Nature Methods* **15**:
535 475–476.
- 536 Guo Y, Li J, Li CI, Long J, Samuels DC, and Shyr Y. 2012. The effect of strand bias in Illumina short-read
537 sequencing data. *BMC Genomics* **13**: 666.
- 538 Iwahara K, Kuriyama T, Shimura S, Williams DW, Yanagisawa M, Nakagawa K, and Karasawa T. 2006.
539 Detection of *cfxA* and *cfxA2*, the β -Lactamase Genes of *Prevotella* spp., in Clinical Samples from Den-
540 toalveolar Infection by Real-Time PCR. *Journal of Clinical Microbiology* **44**: 172–176.
- 541 Jayasundara D, Saeed I, Maheswararajah S, Chang B, Tang SL, and Halgamuge SK. 2015. ViQuaS: An
542 improved reconstruction pipeline for viral quasispecies spectra generated by next-generation sequencing.
543 *Bioinformatics* **31**: 886–896.
- 544 Katoh K, Misawa K, Kuma Ki, and Miyata T. 2002. MAFFT: A novel method for rapid multiple sequence
545 alignment based on fast Fourier transform. *Nucleic Acids Research* **30**: 3059–3066.
- 546 Kazantseva E, Donmez A, Frolova M, Pop M, and Kolmogorov M. 2024. Strainy: Phasing and assembly of
547 strain haplotypes from long-read metagenome sequencing. *Nature Methods* pp. 1–10.
- 548 Kinloch NN, Shahid A, Dong W, Kirkby D, Jones BR, Beelen CJ, MacMillan D, Lee GQ, Mota TM,
549 Sudderuddin H, et al.. 2023. HIV reservoirs are dominated by genetically younger and clonally enriched
550 proviruses. *mBio* **14**: e02417–23.
- 551 Knyazev S, Tsyvina V, Shankar A, Melnyk A, Artyomenko A, Malygina T, Porozov YB, Campbell EM,
552 Switzer WM, Skums P, et al.. 2021. Accurate assembly of minority viral haplotypes from next-generation
553 sequencing through efficient noise reduction. *Nucleic Acids Research* **49**: e102.
- 554 Kosakovsky Pond SL, Posada D, Gravenor MB, Woelk CH, and Frost SD. 2006. GARD: A genetic algorithm
555 for recombination detection. *Bioinformatics* **22**: 3096–3098.
- 556 Köster J. 2016. Rust-Bio: A fast and safe bioinformatics library. *Bioinformatics* **32**: 444–446.
- 557 Lacroix JM and Walker CB. 1996. Detection and prevalence of the tetracycline resistance determinant Tet
558 Q in the microbiota associated with adult periodontitis. *Oral Microbiology and Immunology* **11**: 282–288.
- 559 Lancia G, Bafna V, Istrail S, Lippert R, and Schwartz R. 2001. SNPs Problems, Complexity, and Algorithms.
560 In *Algorithms — ESA 2001* (ed. FM auf der Heide), Lecture Notes in Computer Science, pp. 182–193.
561 Springer, Berlin, Heidelberg.
- 562 Letunic I and Bork P. 2021. Interactive Tree Of Life (iTOL) v5: An online tool for phylogenetic tree display
563 and annotation. *Nucleic Acids Research* **49**: W293–W296.
- 564 Li D, Liu CM, Luo R, Sadakane K, and Lam TW. 2015. MEGAHIT: An ultra-fast single-node solution for
565 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**: 1674–1676.

- 566 Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–3100.
- 567 Li H and Durbin R. 2024. Genome assembly in the telomere-to-telomere era. *Nature Reviews Genetics* **25**:
568 658–670.
- 569 Liao YC, Chen FJ, Chuang MC, Wu HC, Ji WC, Yu GY, and Huang TS. 2022. High-Integrity Sequencing
570 of Spike Gene for SARS-CoV-2 Variant Determination. *International Journal of Molecular Sciences* **23**:
571 3257.
- 572 Liu D and Steinegger M. 2021. Block aligner: Fast and flexible pairwise sequence alignment with SIMD-
573 accelerated adaptive blocks. Preprint, Bioinformatics.
- 574 Luo X, Kang X, and Schönhuth A. 2022. Strainline: Full-length de novo viral haplotype reconstruction from
575 noisy long reads. *Genome Biology* **23**: 29.
- 576 Markov PV, Ghafari M, Beer M, Lythgoe K, Simmonds P, Stilianakis NI, and Katzourakis A. 2023. The
577 evolution of SARS-CoV-2. *Nature Reviews Microbiology* **21**: 361–379.
- 578 McElroy K, Zagordi O, Bull R, Luciani F, and Beerenwinkel N. 2013. Accurate single nucleotide variant
579 detection in viral populations by combining probabilistic clustering with a statistical test of strand bias.
580 *BMC Genomics* **14**: 501.
- 581 Meyer F, Fritz A, Deng ZL, Koslicki D, Lesker TR, Gurevich A, Robertson G, Alser M, Antipov D, Beghini
582 F, et al.. 2022. Critical Assessment of Metagenome Interpretation: The second round of challenges. *Nature*
583 *Methods* **19**: 429–440.
- 584 Meyerhans A, Vartanian JP, and Wain-Hobson S. 1990. DNA recombination during PCR. *Nucleic Acids*
585 *Research* **18**: 1687–1691.
- 586 Mori M, Ode H, Kubota M, Nakata Y, Kasahara T, Shigemi U, Okazaki R, Matsuda M, Matsuoka K,
587 Sugimoto A, et al.. 2022. Nanopore Sequencing for Characterization of HIV-1 Recombinant Forms. *Mi-*
588 *crobiology Spectrum* **10**: e01507–22.
- 589 Myers EW. 2005. The fragment assembly string graph. *Bioinformatics (Oxford, England)* **21 Suppl 2**:
590 ii79–85.
- 591 Nimsamer P, Sawaswong V, Klomkiew P, Kaewsapsak P, Puenpa J, Poovorawan Y, and Payungporn S.
592 2023. “Nano COVID-19”: Nanopore sequencing of spike gene to identify SARS-CoV-2 variants of concern.
593 *Experimental Biology and Medicine* **248**: 1841–1849.
- 594 Olkkola S, Juntunen P, Heiska H, Hyytiäinen H, and Hänninen ML. 2010. Mutations in the rpsL gene are
595 involved in streptomycin resistance in *Campylobacter coli*. *Microbial Drug Resistance (Larchmont, N.Y.)*
596 **16**: 105–110.

- 597 Patterson M, Marschall T, Pisanti N, van Iersel L, Stougie L, Klau GW, and Schönhuth A. 2015. What-
598 sHap: Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *Journal of Computational*
599 *Biology: A Journal of Computational Molecular Cell Biology* **22**: 498–509.
- 600 Price MN, Dehal PS, and Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees for Large
601 Alignments. *PLOS ONE* **5**: e9490.
- 602 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, and Mesirov JP. 2011.
603 Integrative Genomics Viewer. *Nature biotechnology* **29**: 24–26.
- 604 Rognes T, Flouri T, Nichols B, Quince C, and Mahé F. 2016. VSEARCH: A versatile open source tool for
605 metagenomics. *PeerJ* **4**: e2584.
- 606 Ronen R, Boucher C, Chitsaz H, and Pevzner P. 2012. SEQuel: Improving the accuracy of genome assemblies.
607 *Bioinformatics* **28**: i188–i196.
- 608 Rubner Y, Tomasi C, and Guibas L. 1998. A metric for distributions with applications to image databases.
609 In *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)*, pp. 59–66.
- 610 Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA, Wollenberg RD, and Albertsen M. 2022.
611 Oxford Nanopore R10.4 long-read sequencing enables the generation of near-finished bacterial genomes
612 from pure cultures and metagenomes without short-read or reference polishing. *Nature Methods* **19**:
613 823–826.
- 614 Shaw J, Gounot JS, Chen H, Nagarajan N, and Yu YW. 2024. Floria: Fast and accurate strain haplotyping
615 in metagenomes.
- 616 Shaw J and Yu YW. 2022. Flopp: Extremely Fast Long-Read Polyploid Haplotype Phasing by Uniform Tree
617 Partitioning. *Journal of Computational Biology* **29**: 195–211.
- 618 Shaw J and Yu YW. 2023. Fast and robust metagenomic sequence comparison through sparse chaining with
619 skani. *Nature Methods* pp. 1–5.
- 620 Shay JA, Haniford LSE, Cooper A, Carrillo CD, Blais BW, and Lau CHF. 2023. Exploiting a targeted resis-
621 tome sequencing approach in assessing antimicrobial resistance in retail foods. *Environmental Microbiome*
622 **18**: 25.
- 623 Slizovskiy IB, Oliva M, Settle JK, Zyskina LV, Prospero M, Boucher C, and Noyes NR. 2022. Target-enriched
624 long-read sequencing (TELSeq) contextualizes antimicrobial resistance genes in metagenomes. *Microbiome*
625 **10**: 185.
- 626 Stevenson KR, Coolon JD, and Wittkopp PJ. 2013. Sources of bias in measures of allele-specific expression
627 derived from RNA-seq data aligned to a single reference genome. *BMC Genomics* **14**: 536.
- 628 Traag VA, Waltman L, and van Eck NJ. 2019. From Louvain to Leiden: Guaranteeing well-connected
629 communities. *Scientific Reports* **9**: 5233.

- 630 Tyson JR, James P, Stoddart D, Sparks N, Wickenhagen A, Hall G, Choi JH, Lapointe H, Kamelian K,
631 Smith AD, et al.. 2020. Improvements to the ARTIC multiplex PCR method for SARS-CoV-2 genome
632 sequencing using nanopore. *bioRxiv: The Preprint Server for Biology* p. 2020.09.04.283077.
- 633 Van Rossum T, Ferretti P, Maistrenko OM, and Bork P. 2020. Diversity within species: Interpreting strains
634 in microbiomes. *Nature Reviews Microbiology* **18**: 491–506.
- 635 Vedantam G, Guay GG, Austria NE, Doktor SZ, and Nichols BP. 1998. Characterization of mutations
636 contributing to sulfathiazole resistance in *Escherichia coli*. *Antimicrobial Agents and Chemotherapy* **42**:
637 88–93.
- 638 Vicedomini R, Quince C, Darling AE, and Chikhi R. 2021. Strawberry: Automated strain separation in
639 low-complexity metagenomes using long reads. *Nature Communications* **12**: 4485.
- 640 Warburton PJ, Amodeo N, and Roberts AP. 2016. Mosaic tetracycline resistance genes encoding ribosomal
641 protection proteins. *Journal of Antimicrobial Chemotherapy* **71**: 3333–3339.
- 642 Wick RR. 2019. Badread: Simulation of error-prone long reads. *Journal of Open Source Software* **4**: 1316.
- 643 Wick RR, Judd LM, Gorrie CL, and Holt KE. 2017. Completing bacterial genome assemblies with multiplex
644 MinION sequencing. *Microbial Genomics* **3**: e000132.
- 645 Wilm A, Aw PPK, Bertrand D, Yeo GHT, Ong SH, Wong CH, Khor CC, Petric R, Hibberd ML, and
646 Nagarajan N. 2012. LoFreq: A sequence-quality aware, ultra-sensitive variant caller for uncovering cell-
647 population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Research* **40**: 11189–
648 11201.
- 649 Zagordi O, Bhattacharya A, Eriksson N, and Beerenwinkel N. 2011. ShoRAH: Estimating the genetic
650 diversity of a mixed sample from next-generation sequencing data. *BMC Bioinformatics* **12**: 119.
- 651 Zhang T, Li H, Ma S, Cao J, Liao H, Huang Q, and Chen W. 2023. The newest Oxford Nanopore R10.4.1
652 full-length 16S rRNA sequencing enables the accurate resolution of species-level microbial community
653 profiling. *Applied and Environmental Microbiology* **89**: e00605–23.
- 654 Zhou Q, Ji F, Lin D, Liu X, Zhu Z, and Ruan J. 2024. KSNP: A fast de Bruijn graph-based haplotyping
655 tool approaching data-in time cost. *Nature Communications* **15**: 3126.