

# 1 Nucleotide-resolution bacterial pan-genomics with 2 reference graphs

3 Rachel M Colquhoun<sup>1,2,3</sup>, Michael B Hall<sup>1</sup>, Leandro Lima<sup>1</sup>, Leah W Roberts<sup>1</sup>, Kerri M  
4 Malone<sup>1</sup>, Martin Hunt<sup>1,4</sup>, Brice Letcher<sup>1</sup>, Jane Hawkey<sup>5</sup>, Sophie George<sup>4</sup>, Louise Pankhurst<sup>3,6</sup>,  
5 Zamin Iqbal<sup>1,\*</sup>

## 6 Affiliations:

- 7 1. European Bioinformatics Institute, Hinxton, Cambridge CB10 1SD, UK
- 8 2. Wellcome Trust Centre for Human Genetics, University of Oxford, Roosevelt Drive, Oxford, UK
- 9 3. Institute of Evolutionary Biology, Ashworth Laboratories, University of Edinburgh, UK
- 10 4. Nuffield Department of Medicine, University of Oxford, Oxford, UK
- 11 5. Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria  
12 3004, Australia
- 13 6. Department of Zoology, Mansfield Road, University of Oxford

14 \* corresponding author, email [zi@ebi.ac.uk](mailto:zi@ebi.ac.uk)

## 15 Abstract

16 Bacterial genomes follow a U-shaped frequency distribution whereby most genomic loci are  
17 either rare (accessory) or common (core) - the alignable fraction of two genomes from a  
18 single species might be only 50%. Standard tools therefore analyse mutations only in the  
19 core genome, ignoring accessory mutations.

20 We present a novel pan-genome graph structure and algorithms implemented in the  
21 software *pandora*, which approximates a sequenced genome as a recombinant of reference  
22 genomes, detects novel variation and then pan-genotypes multiple samples.

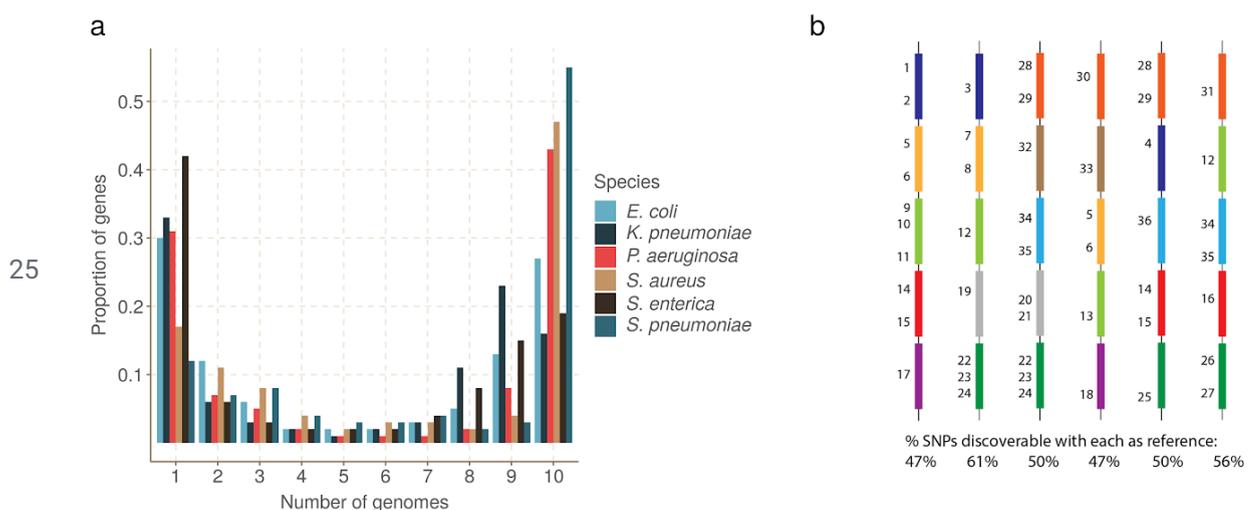
23 Constructing a reference graph from 578 *E. coli* genomes, we analyse a diverse set of 20 *E.*  
24 *coli* isolates. We show, for rare variants, *pandora* recovers at least 13k more SNPs than  
25 single-reference based tools, achieving equal or better error rates with Nanopore as with  
26 Illumina data, and providing a stable framework for analysing diverse samples without  
27 reference bias. This is a significant step towards comprehensive analyses of bacterial  
28 genetic variation.

## 29 Introduction

30 Bacterial genomes evolve by multiple mechanisms including: mutation during replication,  
31 allelic and non-allelic homologous recombination. These processes result in a population of  
32 genomes that are mosaics of each other. Given multiple contemporary genomes, the

1 segregating variation between them allows inferences to be made about their evolutionary  
2 history. These analyses are central to the study of bacterial genomics and evolution<sup>1-4</sup> with  
3 different questions requiring focus on separate aspects of the mosaic: fine-scale (mutations)  
4 or coarse (gene presence, synteny). In this paper, we provide a new and accessible  
5 conceptual model that combines both fine and coarse bacterial variation. Using this new  
6 understanding to better represent variation, we can access previously hidden single  
7 nucleotide polymorphisms (SNPs), insertions and deletions (indels).

8 Genes cover 85-90% of bacterial genomes<sup>5</sup>, and shared gene content is commonly used as  
9 a measure of whole-genome similarity. In fact, the full set of genes present in a species - the  
10 *pan-genome* - is in general much larger than the number found in any single genome. A  
11 frequency distribution plot of genes within a set of bacterial genomes has a characteristic  
12 asymmetric U-shaped curve<sup>6-10</sup>, as shown in Figure 1a. As a result, a collection of  
13 *Escherichia coli* genomes might only have 50% of their genes (and therefore their whole  
14 genome)<sup>3</sup> in common. This highlights a limitation in the standard approach to analysing  
15 genetic variation, whereby a single genome is treated as a reference, and all other genomes  
16 are interpreted as differences from it. In bacteria, a single reference genome will inevitably  
17 lack many of the genes in the pan-genome, and completely miss genetic variation therein  
18 (Figure 1b). We call this *hard reference bias*, to distinguish from the more common concern,  
19 that increased divergence of a reference from the genome under study leads to  
20 read-mapping problems, which we term *soft reference bias*. The standard workaround for  
21 these issues in bacterial genomics is to restrict analysis either to very similar genomes using  
22 a closely related reference (e.g. in an outbreak) or to analyse SNPs only in the core genome  
23 (present in most samples) and outside the core to simply study presence/absence of  
24 genes<sup>11</sup>.



26 **Figure 1. Universal gene frequency distribution in bacteria and the single-reference**  
27 **problem.** a) Frequency distribution of genes in 10 genomes of 6 bacterial species  
28 (*Escherichia coli*, *Klebsiella pneumoniae*, *Pseudomonas aeruginosa*, *Staphylococcus*  
29 *aureus*, *Salmonella enterica* and *Streptococcus pneumoniae*) showing the characteristic  
30 U-shaped curve - most genes are rare or common. b) Illustrative depiction of the  
31 single-reference problem, a consequence of the U-shaped distribution. Each vertical column

1 *is a bacterial genome, and each coloured bar is a gene. Numbers are identifiers for SNPs -*  
2 *there are 50 in total. Thus the dark blue gene has 4 SNPs numbers 1-4. This figure does not*  
3 *detail which genome has which allele. Below each column is the proportion of SNPs that are*  
4 *discoverable when that genome is used as a reference genome. Because no single*  
5 *reference contains all the genes in the population, it can only access a fraction of the SNPs.*

6 In this study we address the variation deficit caused by a single-reference approach. Given  
7 Illumina or Nanopore sequence data from potentially divergent isolates of a bacterial  
8 species, we attempt to detect all of the variants between them. Our approach is to  
9 decompose the pan-genome into atomic units (loci) which tend to be preserved over  
10 evolutionary timescales. Our loci are genes and intergenic regions in this study, but the  
11 method is agnostic to such classifications, and one could add any other grouping wanted  
12 (e.g. operons or mobile genetic elements). Instead of using a single genome as a reference,  
13 we collect a panel of representative reference genomes and use them to construct a set of  
14 reference graphs, one for each locus. Reads are mapped to this set of graphs and from this  
15 we are able to discover and genotype variation. By letting go of prior information on locus  
16 ordering in the reference panel, we are able to recognise and genotype variation in a locus  
17 regardless of its wider context. Since Nanopore reads are typically long enough to  
18 encompass multiple loci, it is possible to subsequently infer the order of loci - although that is  
19 outside the scope of this study.

20 The use of graphs as a generalisation of a linear reference is an active and maturing  
21 field<sup>12-19</sup>. Much recent graph genome work has gone into showing that genome graphs  
22 reduce the impact of soft reference bias on mapping<sup>12</sup>, and on generalising alignment to  
23 graphs<sup>16,20</sup>. However there has not yet been any study (to our knowledge) addressing SNP  
24 analysis across a diverse cohort, including more variants that can fit on any single reference.  
25 In particular, all current graph methods require a reference genome to be provided in  
26 advance to output genetic variants in the standard Variant Call Format (VCF)<sup>21</sup> - thus  
27 immediately inheriting a hard bias when applied to bacteria (see Figure 1b).

28 We have made a number of technical innovations. First, a recursive clustering algorithm that  
29 converts a multiple sequence alignment (MSA) of a locus into a graph. This avoids the  
30 complexity “blowups” that plague graph genome construction from unphased VCF files<sup>12,14</sup>.  
31 Second, a graph representation of genetic variation based on (w,k)-minimizers<sup>22</sup>. Third, using  
32 this representation we avoid unnecessary full alignment to the graph and instead use  
33 quasi-mapping to genotype on the graph. Fourth, discovery of variation missing from the  
34 reference graph using local assembly. Fifth, use of a canonical dataset-dependent reference  
35 genome designed to maximise clarity of description of variants (the value of this will be made  
36 clear in the main text).

37 We describe these below, and evaluate our implementation, *pandora*, on a diverse set of *E.*  
38 *coli* genomes with both Illumina and Nanopore data. We show that, compared with  
39 reference-based approaches, *pandora* recovers a significant proportion of the missing  
40 variation in rare loci, performs much more stably across a diverse dataset, successfully

1 infers a better reference genome for VCF output, and outperforms current tools for Nanopore  
2 data.

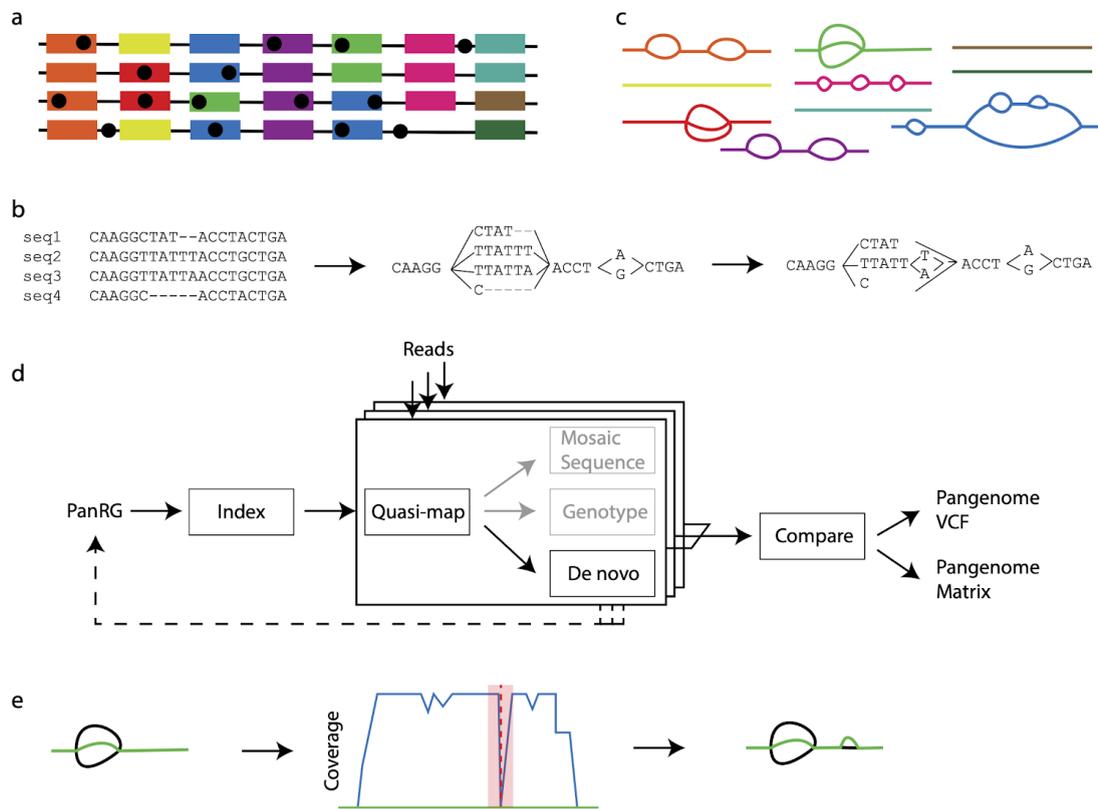
### 3 Results:

#### 4 Pan-genome graph representation

5 We set out to define a generalised reference structure which allows detection of SNPs and  
6 other variants across the whole pan-genome, without attempting to record long-range  
7 structure or coordinates. We define a *Pan-genome Reference Graph* (PanRG) as an  
8 unordered collection of sequence graphs, termed *local graphs*, each of which represents a  
9 locus, such as a gene or intergenic region. Each local graph is constructed from a MSA of  
10 known alleles of this locus, using a recursive cluster-and-collapse (RCC) algorithm  
11 (Supplementary Animation 1: recursive clustering construction). The output is guaranteed to  
12 be a directed acyclic sequence graph allowing hierarchical nesting of genetic variation while  
13 meeting a “balanced parentheses” criterion (see Figure 2b and Methods). Each path through  
14 the graph from source to sink represents a possible recombinant sequence for the locus.  
15 The disjoint nature of this pan-genome reference allows loci such as genes to be compared  
16 regardless of their wider genomic context. We implement this construction algorithm in the  
17 *make\_prg* tool which outputs the graph as a file (see Figures 2a-c, Methods). Subsequent  
18 operations, based on this, are implemented in the software package *pandora*. The overall  
19 workflow is shown in Figure 2.

20 To index a PanRG, we generalise a type of sparse marker k-mer ((w,k)-minimizer),  
21 previously defined for strings, to directed acyclic graphs (see Methods). Each local graph is  
22 *sketched* with minimizing k-mers, and these are then used to construct a new graph (the  
23 k-mer graph) for each local graph from the PanRG. Each minimizing k-mer is a node, and  
24 edges are added between two nodes if they are adjacent minimizers on a path through the  
25 original local graph. This k-mer graph is isomorphic to the original if  $w \leq k$  (and outside the  
26 first and last  $w+k-1$  bases); all subsequent operations are performed on this graph, which, to  
27 avoid unnecessary new terminology, we also call the local graph.

28 A global index maps each minimizing k-mer to a list of all local graphs containing that k-mer  
29 and the positions therein. Long or short reads are approximately mapped (*quasi-mapped*) to  
30 the PanRG by determining the minimizing k-mers in each read. Any of these read  
31 quasi-mappings found in a local graph are called *hits*, and any local graph with sufficient  
32 clustered hits on a read is considered present in the sample.



1 **Figure 2. The *pandora* workflow.** a) reference panel of genomes; colour signifies locus  
 2 (gene or intergenic region) identifier, and blobs are SNPs. b) multiple sequence alignments  
 3 (MSAs) for each locus are made and converted into a directed acyclic graph. c) local graphs  
 4 constructed from the loci in the reference panel. d) Workflow: the collection of local graphs,  
 5 termed the PanRG, is indexed. Reads from each sample under study are independently  
 6 quasi-mapped to the graph, and a determination is made as to which loci are present in  
 7 each sample. In this process, for each locus, a mosaic approximation of the sequence for  
 8 that sample is inferred, and variants are genotyped. e) regions of low coverage are detected,  
 9 and local de novo assembly is used to generate candidate novel alleles missing from the  
 10 graph. Returning to d), the dotted line shows all the candidate alleles from all samples are  
 11 then gathered and added to the MSAs at the start, and the PanRG is updated. Then, reads  
 12 are quasi-mapped one more time, to the augmented PanRG, generating new mosaic  
 13 approximations for all samples and storing coverages across the graphs; no de novo  
 14 assembly is done this time. Finally, all samples are compared, and a VCF file is produced,  
 15 with a per-locus reference that is inferred by *pandora*.

## 16 Initial sequence approximation as a mosaic of references

17 For each locus identified as present in a sample, we initially approximate the sample's  
 18 sequence as a path through the local graph. The result is a mosaic of sequences from the

- 1 reference panel. This path is chosen to have maximal support by reads, using a dynamic
- 2 programming algorithm on the graph induced by its (w,k)-minimizers (details in Methods).
- 3 The result of this process serves as our initial approximation to the genome under analysis.
  
- 4 Improved sequence approximation: modify mosaic by local assembly
- 5 At this point, we have quasi-mapped reads, and approximated the genome by finding the
- 6 closest mosaic in the graph; however, we expect the genome under study to contain variants
- 7 that are not present in the PanRG. Therefore, to allow discovery of novel SNPs and small
- 8 indels that are not in the graph, for each sample and locus we identify regions of the inferred
- 9 mosaic sequence where there is a drop in read coverage (as shown in Figure 2e). Slices of
- 10 overlapping reads are extracted, and a form of *de novo* assembly is performed using a de
- 11 Bruijn graph. Instead of trying to find a single correct path, the de Bruijn graph is traversed
- 12 (see Methods for details) to all feasible candidate novel alleles for the sample. These alleles
- 13 are added to the reference MSA for the locus, and the local graph is updated. If comparing
- 14 multiple samples, the graphs are augmented with all new alleles from all samples at the
- 15 same time.
  
- 16 Optimal VCF-reference construction for multi-genome comparison

17 In the *compare* step of *pandora* (see Figure 2d), we enable continuity of downstream

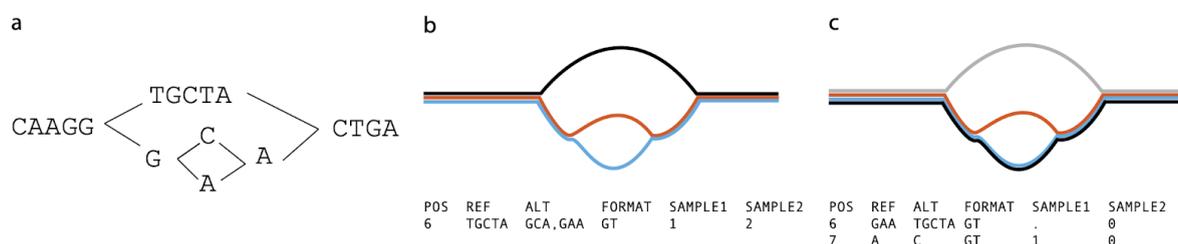
18 analysis by outputting genotype information in the conventional VCF<sup>21</sup>. In this format, each

19 row (record) describes possible alternative allele sequence(s) at a position in a (single)

20 reference genome and information about the type of sequence variant. A column for each

21 sample details the allele seen in that sample, often along with details about the support from

22 the data for each allele.



23 **Figure 3. The representation problem.** a) a local graph. b) The black allele is chosen as

24 reference to enable representation in VCF. The blue/red SNP then requires flanking

25 sequence in order to allow it to have a coordinate. The SNP is thus represented as two ALT

26 alleles, each 3 bases long, and the user is forced to notice they only differ in one base. c)

27 The blue path is chosen as the reference, thus enabling a more succinct and natural

28 representation of the SNP.

1 To output graph variation, we first select a path through the graph to be the reference  
2 sequence and describe any variation within the graph with respect to this path as shown in  
3 Figure 3. We use the chromosome field to detail the local graph within the PanRG in which a  
4 variant lies, and the position field to give the position in the chosen reference path sequence  
5 for that graph. In addition, we output the reference path sequences used as a separate file.

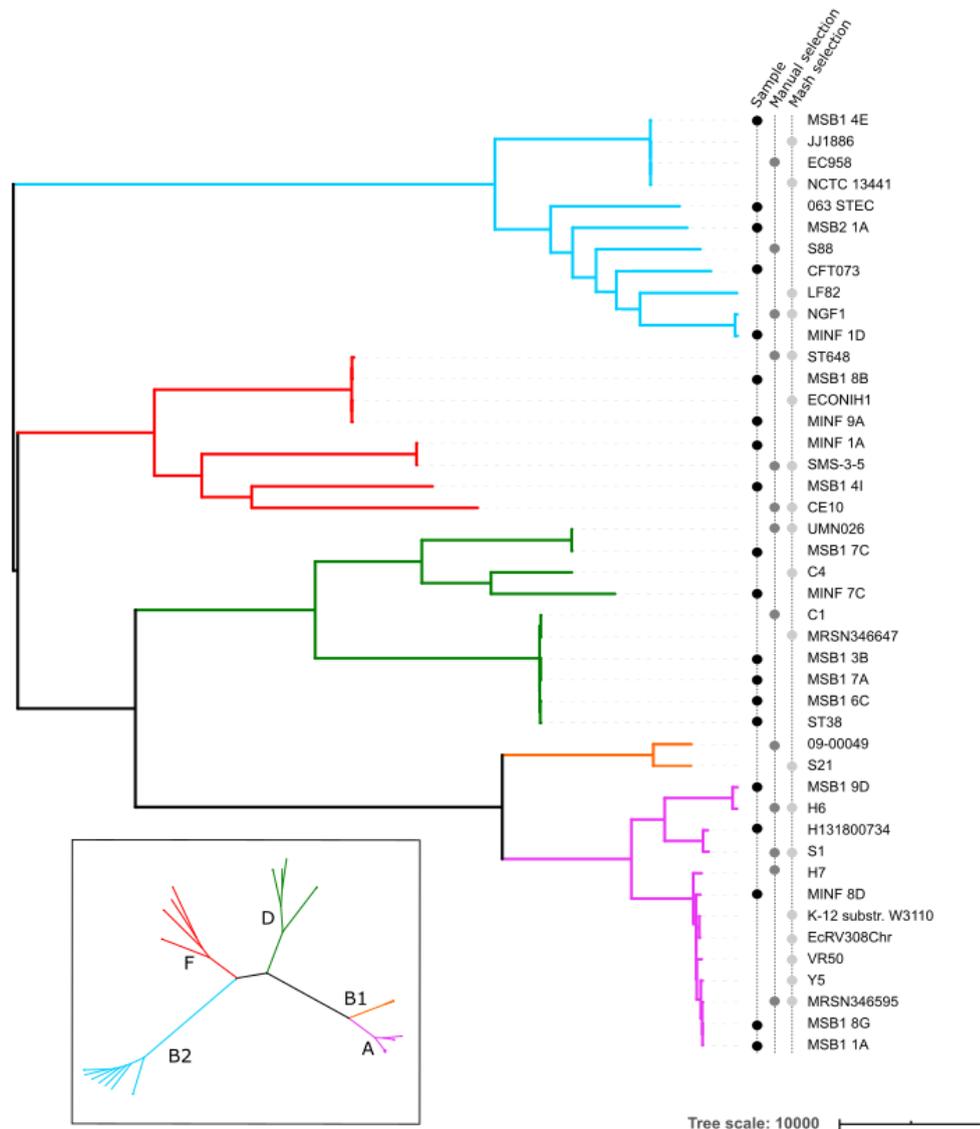
6 For a collection of samples, we want small differences between samples to be recorded as  
7 short alleles in the VCF file rather than longer alleles with shared flanking sequence as  
8 shown in Figure 3b. We therefore choose the reference path for each local graph to be  
9 maximally close to the sample mosaic paths. To do this, we make a copy of the k-mer graph  
10 and increment the coverage along each sample mosaic path, producing a graph with higher  
11 weights on paths shared by more samples. We reuse the mosaic path-finding algorithm (see  
12 Methods) with a modified probability function defined such that the probability of a node is  
13 proportional to the number of samples covering it. This produces a dataset-dependent VCF  
14 reference able to succinctly describe segregating variation in the cohort of genomes under  
15 analysis.

## 16 Constructing a PanRG of *E. coli*

17 We chose to evaluate *pandora* on the recombining bacterial species, *E. coli*, whose  
18 pan-genome has been heavily studied<sup>7,23-26</sup>. MSAs for gene clusters curated with PanX<sup>27</sup>  
19 from 350 RefSeq assemblies were downloaded from <http://pangenome.de> on 3rd May 2018.  
20 MSAs for intergenic region clusters based on 228 *E. coli* ST131 genome sequences were  
21 previously generated with Piggy<sup>28</sup> for their publication. Whilst this panel of intergenic  
22 sequences does not reflect the full diversity within *E. coli*, we included them as an initial  
23 starting point. This resulted in an *E. coli* PanRG containing local graphs for 23,054 genes  
24 and 14,374 intergenic regions. *Pandora* took 24.4h in CPU time (2.3h in runtime with 16  
25 threads) and 12.6 GB of RAM to index the PanRG. As one would expect from the U-shaped  
26 gene frequency distribution, many of the genes were rare in the 578 (=350+228) input  
27 genomes, and so 59%/44% of the genic/intergenic graphs were linear, with just a single  
28 allele.

## 29 Constructing an evaluation set of diverse genomes

30 We first demonstrate that using a PanRG reduces hard bias when comparing a diverse set  
31 of 20 *E. coli* samples by comparison with standard single reference variant callers. We  
32 selected samples from across the phylogeny (including phylogroups A, B2, D and F<sup>29</sup>) where  
33 we were able to obtain both long and short read sequence data from the same isolate.



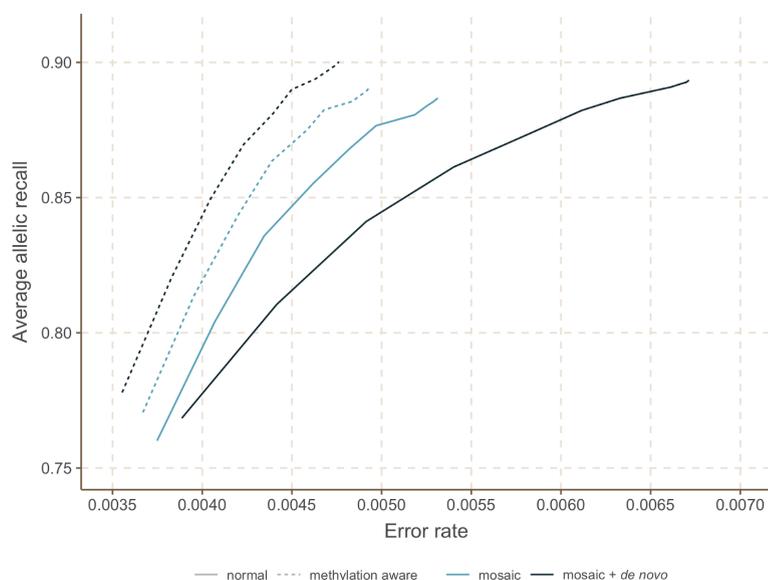
1 **Figure 4. Phylogeny of 20 diverse *E. coli* along with references used for benchmarking**  
2 **single-reference variant callers.** *The 20 E. coli under study are labelled as samples in the*  
3 *left-hand of three vertical label-lines. Phylogroups (clades) are labelled by colour of branch,*  
4 *with the key in the inset. References were selected from RefSeq as being the closest to one*  
5 *of the 20 samples as measured by Mash, or manually selected from a tree (see Methods).*  
6 *Two assemblies from phylogroup B1 are in the set of references, despite there being no*  
7 *sample in that phylogroup.*

8 We used Illumina-polished long read assemblies as truth data, masking positions where the  
9 Illumina data did not support the assembly (see Methods). As comparators, we used  
10 SAMtools<sup>30</sup> (the “classical” variant-caller based on pileups) and Freebayes<sup>31</sup> (a  
11 haplotype-based caller which reduces soft reference bias, wrapped by Snippy  
12 <https://github.com/tseemann/snippy>) for Illumina data, and Medaka (unpublished,  
13 <https://github.com/Nanoporetech/medaka>) and Nanopolish<sup>32</sup> for Nanopore data. In all cases,  
14 we ran the reference-based callers with 24 carefully selected reference genomes (see

1 Methods, and Figure 4). We defined a “truth set” of 618,305 segregating variants by  
2 performing all pairwise whole genome alignments of the 20 truth assemblies, collecting SNP  
3 variants between the pairs, and deduplicating them by clustering into equivalence classes.  
4 Each class, or *pan-variant*, represents the same variant found at different coordinates in  
5 different genomes (see Methods). We evaluated error rate, pan-variant recall (PVR,  
6 proportion of truth set discovered) and average allelic recall (AvgAR, average of the  
7 proportion of alleles of each pan-variant that are found). To clarify the definitions, consider a  
8 toy example. Suppose we have three genes, each with one SNP between them. The first  
9 gene is rare, present in 2/20 genomes. The second gene is at an intermediate frequency, in  
10 10/20 genomes. The third is a strict core gene, present in all genomes. The SNP in the first  
11 gene has alleles A,C at 50% frequency (1 A and 1 C). The SNP in the second gene has  
12 alleles G,T at 50% frequency (5 G and 5 T). The SNP in the third gene has alleles A,T with  
13 15 A and 5 T. Suppose a variant caller found the SNP in the first gene, detecting the two  
14 correct alleles. For the second gene’s SNP, it detected only one G and one T, failing to  
15 detect either allele in the other 8 genomes. For the third gene’s SNP, it detected all the 5 T’s,  
16 but no A. Here, the pan-variant recall would be:  $(1 + 1 + 0) / 3 = 0.66$  - *i.e.* score a 1 if both  
17 alleles are found, irrespective of how often- and the average allelic recall would be  $(2/2 +$   
18  $2/10 + 5/20)/3=0.48$ .

## 19 Methylation-aware basecalling improves results

20 In Figure 5, we show for 4 samples the effect of methylation-aware Nanopore basecalling on  
21 the AvgAR/error rate curve for *pandora* with/without novel variant discovery via local  
22 assembly.

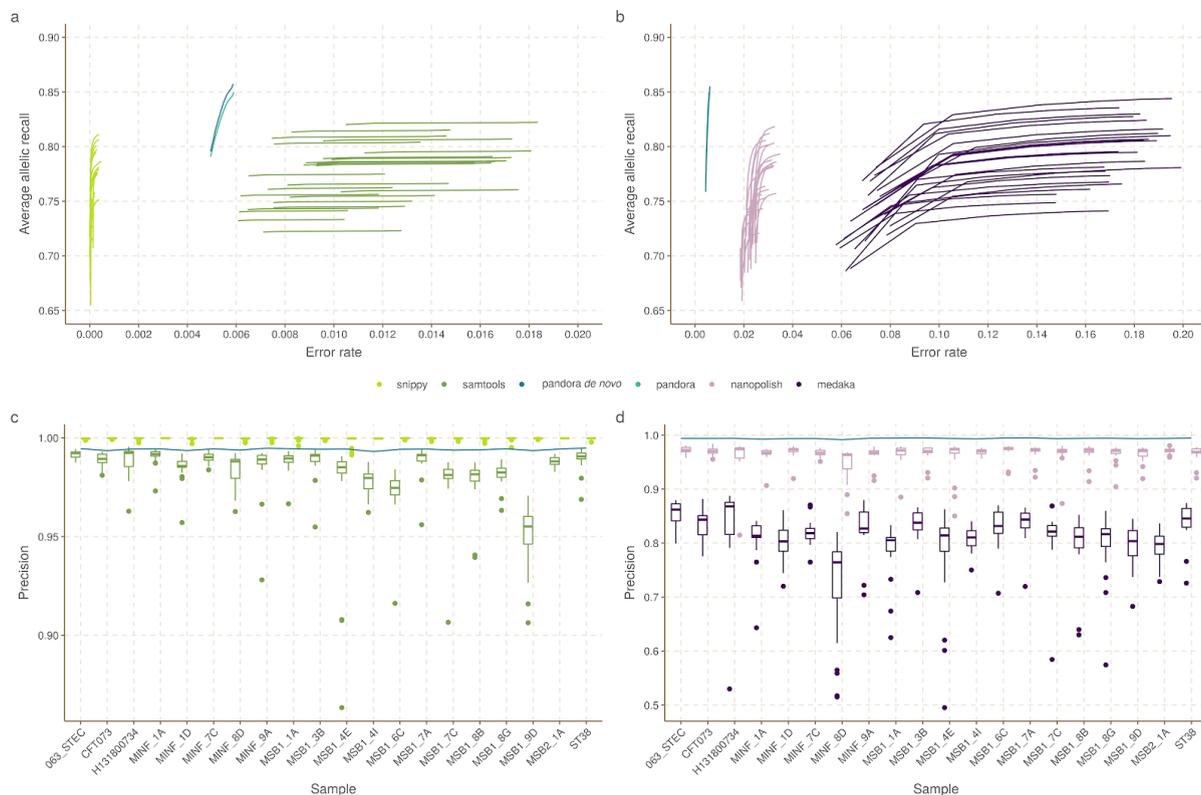


23 **Figure 5. The effect of methylation-aware basecalling on local *de novo* assembly.** We  
24 show the Average Allelic Recall and Error Rate curve for *pandora* with normal (solid line) or  
25 methylation-aware (dashed line) Guppy basecalling on 4 out of the 20 samples. For each of  
26 these input data, we show results for Pandora’s first approximation to a genome as a mosaic  
27 (recombinant) of the input reference panel (mosaic, light blue), and then the improved  
28 approximation with added *de novo* discovery (mosaic+*de novo*, dark blue).

1 The top right of each curve corresponds to completely unfiltered results; increasing the  
 2 genotype confidence threshold (see Methods) moves each curve towards the bottom-left,  
 3 increasing precision at the cost of recall. Notably, with normal basecalling, local *de novo*  
 4 assembly increases the error rate from 0.53% to 0.67%, with a negligible increase in recall,  
 5 from 88.7% to 89.3%, whereas with methylation-aware basecalling it increases the recall  
 6 from 89.1% to 90% and slightly decreases the error rate from 0.49% to 0.48%. On the basis  
 7 of this, from here on we work entirely with reads that are basecalled with a  
 8 methylation-aware model, and move to the full dataset of 20 samples.

## 9 Benchmarking recall, error rate and dependence on reference

10 We show in Figures 6a,b the Illumina and Nanopore AvgAR/recall plots for *pandora* and four  
 11 single-reference tools with no filters applied. For all of these, we modify only the minimum  
 12 genotype confidence to move up and down the curves (see Methods).



13 **Figure 6. Benchmarks of recall/error and dependence of precision on reference**  
 14 **genome, for *pandora* and other tools on 20-way dataset.** a) The average allelic recall  
 15 and error rate curve for *pandora*, *SAMtools* and *snippy* on 100x of Illumina data.  
 16 *Snippy/SAMtools* both run 24 times with the different reference genomes shown in figure 4,  
 17 resulting in multiple lines for each tool (one for each reference) b) The average allelic recall  
 18 and error rate curve for *pandora*, *medaka* and *nanopolish* on 100x of Nanopore data;  
 19 multiple lines for *medaka/nanopolish*, one for each reference genome. Note panels a and b  
 20 have the same y axis scale and limits, but different x axes; c) The precision of *pandora*,  
 21 *SAMtools* and *snippy* on 100x of Illumina data. The boxplots show the distribution of

1 *SAMtools*' and *snippy*'s precision depending on which of the 24 references was used, and  
2 the blue line connects *pandora*'s results; d) The precision of *pandora* (line plot), *medaka* and  
3 *nanopolish* (both boxplots) on 100x of Nanopore data. Note different y axis scale/limits in  
4 panels c,d.

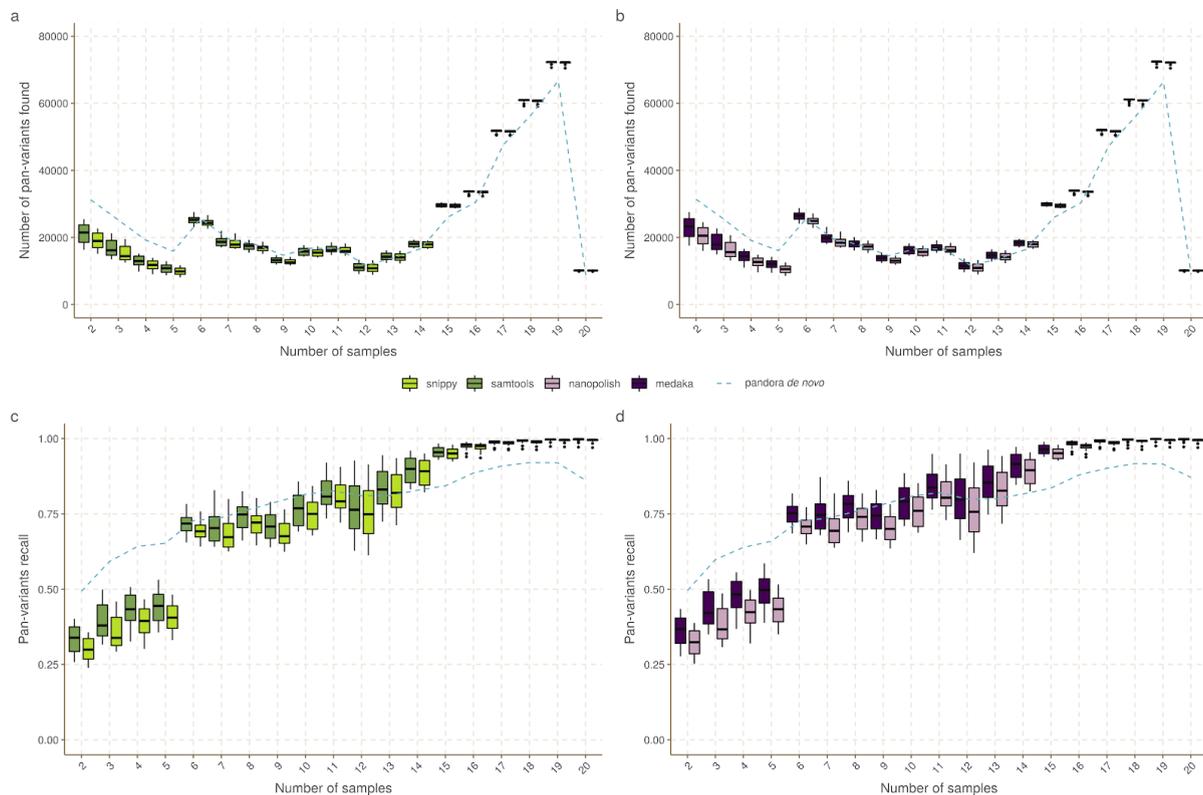
5 We highlight three observations. Firstly, *pandora* achieves essentially the same recall and  
6 error rate for the Illumina and Nanopore data (85% AvgAR and 0.6% error rate at the  
7 top-right of the curve, completely unfiltered). Second, choice of reference has a significant  
8 effect on both AvgAR and error rate for the single-reference callers; the reference which  
9 enables the highest recall does not lead to the best error rate (for *SAMtools* and *medaka* in  
10 particular). Third, *pandora* achieves better AvgAR (86%) than all other tools (all between  
11 81% and 84%, see Supplementary Table 2), and a better error rate (0.6%) than *SAMtools*  
12 (1%), *nanopolish* (2.4%) and *medaka* (14.8%). However, *snippy* achieves a significantly  
13 better error rate than all other tools (0.01%). We confirmed that adding further filters slightly  
14 improved error rates, but did not change the overall picture (Supplementary Figure 1,  
15 Methods, Supplementary Table 2). The results are also in broad agreement if the PVR is  
16 plotted instead of AvgAR (Supplementary Figure 2). However, these AvgAR and PVR figures  
17 are hard to interpret because *pandora* and the reference-based tools have recall that varies  
18 differently across the locus frequency spectrum - we explore this further below.

19 We ascribe the similarity between the Nanopore and Illumina performance of *pandora* to  
20 three reasons. First, the PanRG is a strong prior - our first approximation does not contain  
21 any Nanopore sequence, but simply uses quasi-mapped reads to find the nearest mosaic in  
22 the graph. Second, mapping long Nanopore reads which completely cover entire genes is  
23 easier than mapping Illumina data, and allows us to filter out erroneous k-mers within reads  
24 after deciding when a gene is present. Third, this performance is only achieved when we use  
25 methylation-aware basecalling of Nanopore reads, presumably removing most systematic  
26 bias (see Figure 5).

27 In Figure 6c,d we show for Illumina and Nanopore data, the impact of reference choice on  
28 the precision of calls on each of the 20 samples. While precision is consistent across all  
29 samples for *pandora*, we see a dramatic effect of reference-choice on precision of *SAMtools*,  
30 *medaka* and *nanopolish*. The effect is also detectable for *snippy*, but to a much lesser  
31 extent.

32 Finally, we measured the performance of locus presence detection, restricting to  
33 genes/intergenic regions in the PanRG, so that in principle perfect recall would be possible  
34 (see Methods). In Supplementary Figure 3 we show the distribution of locus presence calls  
35 by *pandora*, split by length of locus for Illumina and Nanopore data. Overall, 93.8%/94.3% of  
36 loci were correctly classified as present or absent for Illumina/Nanopore respectively.  
37 Misclassifications were concentrated on small loci (below 500bp). While 59.2%/57.4% of all  
38 loci in the PanRG are small, 75.5%/74.8% of false positive calls and 98.7%/98.1% of false  
39 negative calls are small loci (see Supplementary Figure 3).

- 1 *Pandora* detects rare variation inaccessible to single-reference methods
- 2 Next, we evaluate the key deliverable of *pandora* - the ability to access genetic variation
- 3 within the accessory genome. We plot this in Figure 7, showing PVR of SNPs in the truth set
- 4 which overlap genes or intergenic regions from the PanRG, broken down by the number of
- 5 samples the locus is present in.

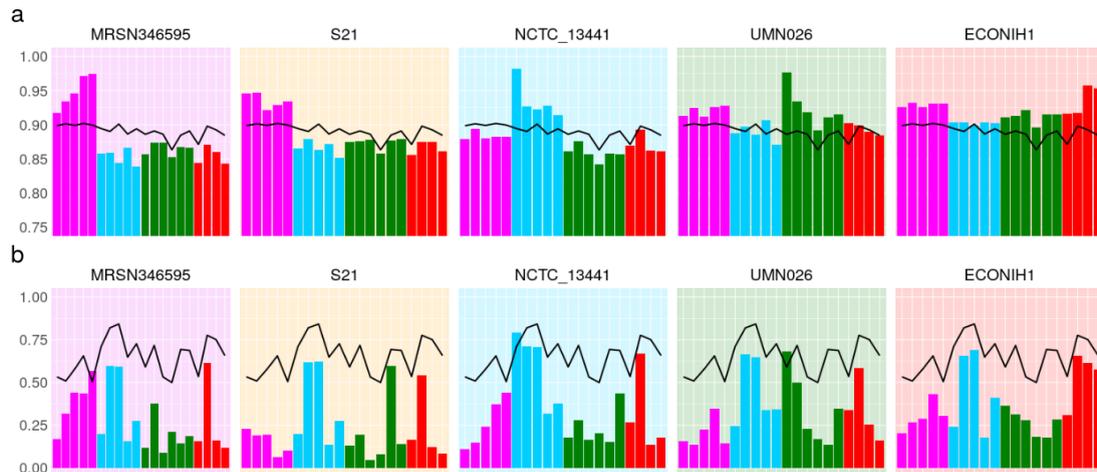


6 **Figure 7. Pan-variant recall across the locus frequency spectrum. Every SNP occurs in**  
 7 **a locus, which is present in some subset of the full set of 20 genomes. In all panels the**  
 8 **SNPs in the golden truth set are broken down by the number of samples the locus is present**  
 9 **in. Left panels (a, c) show results for *pandora* (dotted line), *snippy* and *SAMtools* with**  
 10 ***Illumina* data. Right panels (b, d) show results for *pandora*, *nanopolish* and *medaka* with**  
 11 ***Nanopore* data. Top panels (a, b) show the absolute count of pan-variants found; Bottom**  
 12 **panels (c, d) show the proportion of pan-variants found.**

13 If we restrict our attention to rare variants (present only in 2-5 genomes), we find *pandora*  
 14 recovers at least 19644/26674/13108/22331 more SNPs than  
 15 *SAMtools/snippy/medaka/nanopolish* respectively. As a proportion of rare SNPs in the truth  
 16 set, this is a lift in PVR of 12/17/8/14% respectively. If, instead of pan-variant recall, we look  
 17 at the variation of AvgAR across the locus frequency spectrum (see Supplementary Figure  
 18 4), the gap between *pandora* and the other tools on rare loci is even larger. These  
 19 observations, and Figure 6, confirm and quantify the extent to which we are able to recover  
 20 accessory genetic variation that is inaccessible to single-reference based methods.

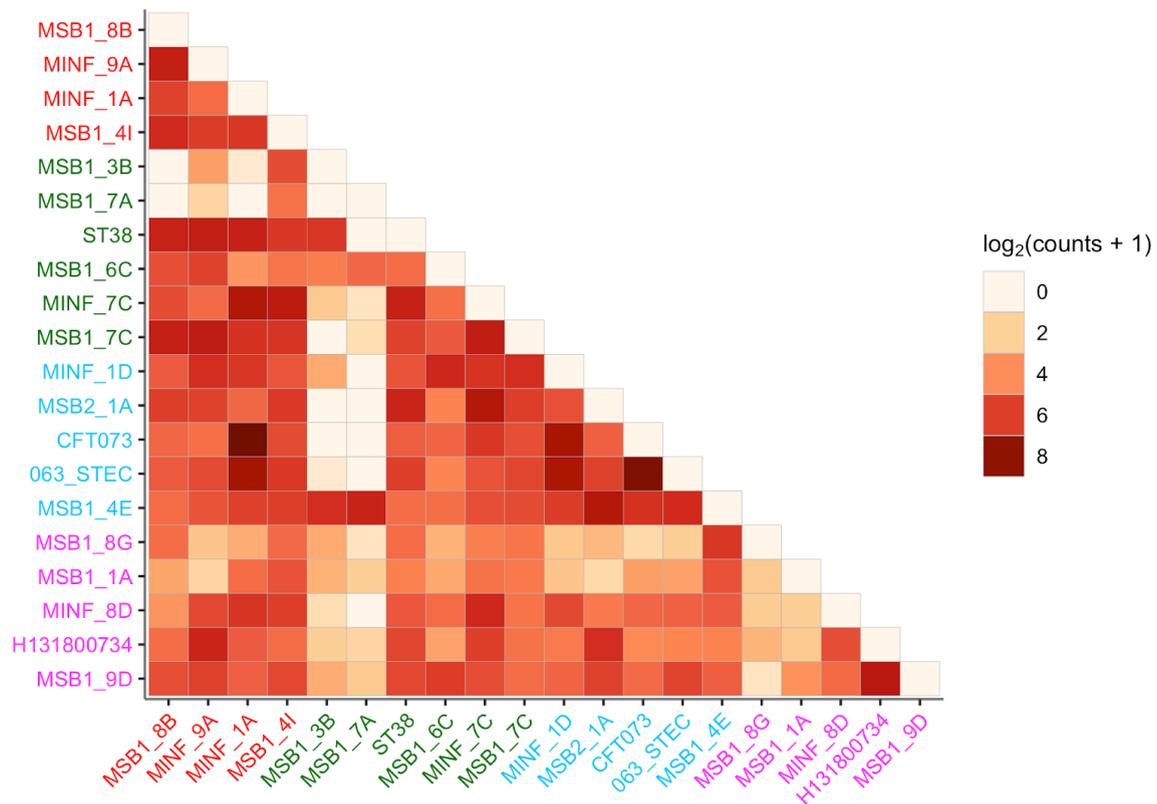
## 1 *Pandora* has consistent results across *E. coli* phylogroups

- 2 We measure the impact of reference bias (and population structure) by quantifying how
- 3 recall varies in phylogroups A, B2, D, and F dependent on whether the reference genome
- 4 comes from the same phylogroup.



5 **Figure 8. Single reference callers achieve higher recall for samples in the same**  
6 **phylogroup as the reference genome, but not for rare loci. a) *pandora* recall (black line)**  
7 **and snippy recall (coloured bars) on the 20 samples; each histogram corresponds to the use**  
8 **of one of 5 exemplar references, one from each phylogroup. The background colour denotes**  
9 **the reference's phylogroup (see Figure 4 inset); note that phylogroup B1 (yellow**  
10 **background) is an outgroup, containing no samples in this dataset; b) Same as a) but**  
11 **restricted to SNPs present in precisely two samples (i.e. where 18 samples have neither**  
12 **allele because the entire locus is missing). Note the differing y-axis limits in the two panels.**

13 We plot the results for *snippy* with 5 exemplar references in Figure 8a (results for all tools  
14 and for all references are in Supplementary Figures 5-8), showing that single references give  
15 5-10% higher recall for samples in their own phylogroup than other phylogroups. By  
16 comparison, *pandora*'s recall is much more consistent, staying stable at ~89% for all  
17 samples regardless of phylogroup. References in phylogroups A and B2 achieve higher  
18 recall in their own phylogroup, but consistently worse than *pandora* for samples in the other  
19 phylogroups (in which the reference does not lie). References in the external phylogroup B1,  
20 for which we had no samples in our dataset, achieve higher recall for samples in the nearby  
21 phylogroup A (see inset, Figure 4), but lower than *pandora* for all others. We also see that  
22 choosing a reference genome from phylogroup F (red), which sits intermediate to the other  
23 phylogroups, provides the most uniform recall across other groups - 2-5% higher than  
24 *pandora*.



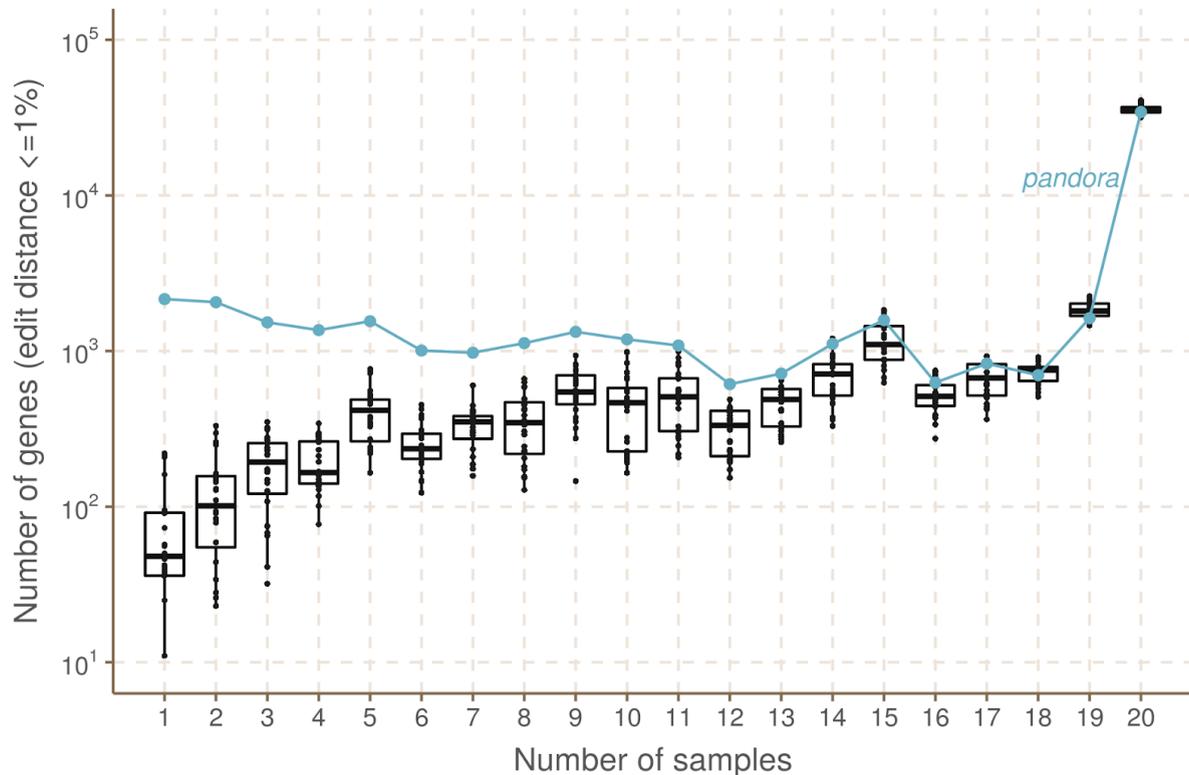
1 **Figure 9. Sharing of variants present in precisely 2 genomes, showing which pairs of**  
2 **genomes they lie in and which phylogroups; darker colours signify higher counts (log**  
3 **scale). Genomes are coloured by their phylogroup (see Figure 4 inset).**

4 These results will, however, be dominated by the shared, core genome. If we replot Figure  
5 8a, restricting to variants in loci present in precisely 2 genomes (abbreviated to 2-variants;  
6 Figure 8b), we find that *pandora* achieves 50-84% recall for each sample (complete data in  
7 Supplementary Figure 9). By contrast, for any choice of reference genome, the results for  
8 single-reference callers vary dramatically per sample. Most samples have recall under 25%,  
9 and there is no pattern of improved recall for samples in the same phylogroup as the  
10 reference. Following up that last observation, if we look at which pairs of genomes share  
11 2-variants (Figure 9), we find there is no enrichment within phylogroups at all. This simply  
12 confirms in our data that presence of rare loci is not correlated with the overall phylogeny.

### 13 *Pandora* VCF reference is closer to samples than any single reference

14 The relationship between phylogenetic distance and gene repertoire similarity is not linear. In  
15 fact, 2 genomes in different phylogroups may have more similar accessory genes than 2 in  
16 the same phylogroup - as illustrated in the previous section (also see Figure 3 in Rocha<sup>3</sup>). As  
17 a result, it is unclear *a priori* how to choose a good reference genome for comparison of  
18 accessory loci between samples. *Pandora* specifically aims to construct an appropriate  
19 reference for maximum clarity in VCF representation. We evaluate how well *pandora* is able

- 1 to find a VCF reference close to the samples under study as follows. We first identified the
- 2 location of all loci in all the 20 sample assemblies and the 24 references (see Methods).



3 **Figure 10. How often do references closely approximate a sample?** *pandora* aims to  
4 infer a reference for use in its VCF, which is as close as possible to all samples. We evaluate  
5 the success of this here. The x-axis shows the number of genomes in which a locus occurs.  
6 The y-axis shows the (log-scaled) count of loci in the 20 samples that are within 1% edit  
7 distance (scaled by locus length) of each reference - box plots for the reference genomes,  
8 and line plot for the VCF reference inferred by *pandora*.

9 We then measured the edit distance between each locus in each of the references and the  
10 corresponding version in the 20 samples. We found that the *pandora*'s VCF-reference lies  
11 within 1% edit distance (scaled by locus length) of the sample far more than any of the  
12 references for loci present in  $\leq 14$  samples (Figure 10; note the log scale). The  
13 improvement is much reduced in the core genome; essentially, in the core, a  
14 phylogenetically close reference provides a good approximation, but it is hard to choose a  
15 single reference that provides a close approximation to all rare loci. By contrast, *pandora* is  
16 able to leverage its reference panel, and the dataset under study, to find a good  
17 approximation.

## 1 Computational performance

2 Performance measurements for single-sample analysis by *pandora* and benchmarked tools  
3 are shown in Supplementary Table 3. In short, *pandora* took 3-4 hours per sample (using 16  
4 cores and up to 10.7 GB of RAM), which was slower than *snippy* (0.1h, 4 cores), *SAMtools*  
5 (0.3h, 1 core) and *medaka* (0.3h, 4 cores), but faster than *nanopolish* (4.6h, 16 cores).

6 *Pandora* alone can do joint analysis of multiple samples and this is currently the most  
7 expensive *pandora* step. Parallelising by gene on a compute cluster, it took 8 hours to  
8 augment the PanRG with novel alleles. This was dominated by the Python implementation of  
9 the RCC clustering algorithm (see Methods) and the use of Clustal Omega<sup>33</sup> for MSA. 90%  
10 of loci required less than 30 minutes to process, and the remainder took less than 2 hours  
11 (see Methods). We discuss below how this could be improved. Finally, it took 28/46 hours to  
12 compare the samples (produce the joint VCF file) for Illumina/Nanopore. Mapping comprised  
13 ~10% of the Illumina time, and ~50% of the Nanopore time. Dynamic programming and  
14 genotyping the VCF file took ~90% of the Illumina time, and ~50% of the Nanopore time.

## 15 Discussion

16 Bacteria are the most diverse and abundant cellular life form<sup>34</sup>. Some species are exquisitely  
17 tuned to a particular niche (e.g. obligate pathogens of a single host) while others are able to  
18 live in a wide range of environments (e.g. *E. coli* can live on plants, in the earth, or  
19 commensally in the gut of various hosts). Broadly speaking, a wider range of environments  
20 correlates with a larger pan-genome, and some parts of the gene repertoire are associated  
21 with specific niches<sup>35</sup>. Our perception of a pan-genome therefore depends on our sampling  
22 of the unknown underlying population structure, and similarly the effectiveness of a PanRG  
23 will depend on the choice of reference panel from which it is built.

24 Many examples from different species have shown that bacteria are able to leverage this  
25 genomic flexibility, adapting to circumstance sometimes by using or losing novel genes  
26 acquired horizontally, and at other times by mutation. There are many situations where  
27 precise nucleotide-level variants matter in interpreting pan-genomes. Some examples  
28 include: compensatory mutations in the chromosome reducing the fitness burden of new  
29 plasmids<sup>36-38</sup>; lineage-specific accessory genes with SNP mutations which distinguish  
30 carriage from infection<sup>39</sup>; SNPs within accessory drug resistance genes leading to significant  
31 differences in antibiograms<sup>40</sup>; and changes in CRISPR spacer arrays showing immediate  
32 response to infection<sup>41,42</sup>. However, up until now there has been no automated way of  
33 studying non-core gene SNPs at all; still less a way of integrating them with gene  
34 presence/absence information. *Pandora* solves these problems, allowing detection and  
35 genotyping of core and accessory variants. It also addresses the problem of what reference  
36 to use as a coordinate system, inferring a mosaic “VCF reference” which is as close as  
37 possible to all samples under study. We find this gives more consistent SNP-calling than any  
38 single reference in our diverse dataset. We focussed primarily on Nanopore data when  
39 designing *pandora*, and show it is possible to achieve higher quality SNP calling with this  
40 data than with current Nanopore tools. Together, these results open the door for empirical

1 studies of the accessory genome, and for new population genetic models of the pan-genome  
2 from the perspective of both SNPs and gene gain/loss.

3 Prior graph genome work, focussing on soft reference bias (in humans), has evaluated  
4 different approaches for selecting alleles for addition to a population graph, based on  
5 frequency, avoiding creating new repeats, and avoiding exponential blowup of haplotypes in  
6 clusters of variants<sup>43</sup>. This approach makes sense when you have unphased diploid VCF  
7 files and are considering all recombinants of clustered SNPs as possible. However, this is  
8 effectively saying we consider the recombination rate to be high enough that all  
9 recombinants are possible. Our approach, building from local MSAs and only collapsing  
10 haplotypes when they agree for a fixed number of bases, preserves more haplotype  
11 structure and avoids combinatorial explosion. Another alternative approach was recently  
12 taken by Norri *et al.*<sup>44</sup>, inferring a set of pseudo founder genomes from which to build the  
13 graph.

14 Another issue is how to select the reference panel of genomes in order to minimize hard  
15 reference bias. One cannot escape the U-shaped frequency distribution; whatever reference  
16 panel is chosen, future genomes under study will contain rare genes not present in the  
17 PanRG. Given the known strong population structure in bacteria, and the association of  
18 accessory repertoires with lifestyle and environment, we would advocate sampling by  
19 geography, host species (if appropriate), lifestyle (e.g. pathogenic versus commensal) and/or  
20 environment. In this study we built our PanRG from a biased dataset (RefSeq) which does  
21 not attempt to achieve balance across phylogeny or ecology, limiting our pan-variant recall to  
22 49% for rare variants (see Figure 7c,d). A larger, carefully curated input panel, such as that  
23 from Horesh *et al.*<sup>45</sup>, would provide a better foundation and potentially improve results.

24 A natural question is then to ask if the PanRG should continually grow, absorbing all variants  
25 ever encountered. From our perspective, the answer is no - a PanRG with variants at all  
26 non-lethal positions would be potentially intractable. The goal is not to have every possible  
27 allele in the PanRG - no more than a dictionary is required to contain absolutely every word  
28 that has ever been said in a language. As with dictionaries, there is a trade-off between  
29 completeness and utility, and in the case of bacteria, the language is far richer than English.  
30 The perfect PanRG contains the vast majority of the genes and intergenic regions you are  
31 likely to meet, and just enough breadth of allelic diversity to ensure reads map without too  
32 many mismatches. Missing alleles should be discoverable by local assembly and added to  
33 the graph, allowing multi-sample comparison of the cohort under study. This allows one to  
34 keep the main PanRG lightweight enough for rapid and easy use.

35 We finish with three potential applications of *pandora*. First, the PanRG should provide a  
36 more interpretable substrate for pan-genome-wide Genome-Wide Association Studies, as  
37 current methods are forced to either ignore the accessory genome or reduce it to k-mers or  
38 unitigs<sup>46-48</sup>. Second, if performing prospective surveillance of microbial isolates taken in a  
39 hospital, the PanRG provides a consistent and unchanging reference, which will cope with  
40 the diversity of strains seen without requiring the user to keep switching reference genome.  
41 In a sense it behaves similarly to whole-genome Multi-Locus Sequence Typing (wgMLST)<sup>49</sup>,

1 with more flexibility, support for intergenic regions, and without the all-or-nothing behaviour  
2 when alleles have a novel SNP. Third, if studying a fixed dataset very carefully, then one  
3 may not want to use a population PanRG, as it necessarily will miss some rare accessory  
4 genes in the dataset. In these circumstances, one could construct a reference graph purely  
5 of the genes/intergenic regions present in this dataset.

6 There are a number of limitations to this study. Firstly, *pandora* is not yet a fully-fledged  
7 production tool. There are two steps that constitute bottlenecks in terms of RAM and speed.  
8 The RCC algorithm used for local graph construction is currently implemented in Python.  
9 However, the underlying algorithm is amenable to a much higher performance  
10 implementation, which is now in progress. Also, we use Clustal Omega<sup>33</sup> for the MSA stage,  
11 and there are faster options which we could use, including options for augmenting an MSA  
12 without a complete rebuild (e.g. MAFFT), which is exactly what we need after local assembly  
13 discovers novel alleles. Secondly, we do not see any fundamental reason why the *pandora*  
14 error rate should be worse than Snippy on Illumina data (see Figure 6C), and will be working  
15 to improve this. Finally, by working in terms of atomic loci instead of a monolithic  
16 genome-wide graph, *pandora* opens up graph-based approaches to structurally diverse  
17 species (and eases parallelisation) but at the cost of losing genome-wide ordering. At  
18 present, ordering can be resolved by (manually) mapping *pandora*-discovered genes onto  
19 whole genome assemblies. However the design of *pandora* also allows for gene-ordering  
20 inference: when Nanopore reads cover multiple genes, the linkage between them is stored in  
21 a secondary de Bruijn graph where the alphabet consists of gene identifiers. This results in a  
22 huge alphabet, but the k-mers are almost always unique, dramatically simplifying “assembly”  
23 compared with normal DNA de Bruijn graphs. This work is still in progress and the subject of  
24 a future study. In the meantime, *pandora* provides new ways to access previously hidden  
25 variation.

## 26 Methods

### 27 Local graph construction

28 We construct each local graph in the PanRG from an MSA using an iterative partitioning  
29 process. The resulting sequence graph contains nested bubbles representing alternative  
30 alleles.

31 Let  $A$  be an MSA of length  $n$ . For each row of the MSA  $a = \{a_0, \dots, a_{n-1}\} \in A$  let  
32  $a_{i,j} = \{a_i, \dots, a_{j-1}\}$  be the subsequence of  $a$  in interval  $[i,j)$ . Let  $s(a)$  be the DNA sequence  
33 obtained by removing all non-AGCT symbols. We can partition alignment  $A$  either *vertically*  
34 by partitioning the interval  $[0,n)$  or *horizontally* by partitioning the set of rows of  $A$ . In both  
35 cases, the partition induces a number of sub-alignments.

1 For vertical partitions, we define  $slice_A(i,j) = \{a_{ij} : a \in A\}$ . We say that interval  $[i,j]$  is a  
 2 *match* interval if  $j - i \geq m$ , where  $m = 7$  is the default minimum match length, and there is a  
 3 single non-trivial sequence in the slice, i.e.  $|\{s(a) : a \in slice_A(i,j) \text{ and } s(a) \neq ""\}| = 1$ .  
 4 Otherwise, we call it a *non-match* interval.

5 For horizontal partitions, we use  $K$ -means clustering<sup>50</sup> to divide sequences into increasing  
 6 numbers of clusters  $K = 2, 3, \dots$  until the *inertia*, a measure of the within-cluster diversity, is  
 7 half that of the original full set of sequences. More formally, let  $U$  be the set of all  $m$ -mers  
 8 (substrings of length  $m$ , the minimum match length) in  $\{s(a) : a \in A\}$ . For  $a \in A$  we  
 9 transform sequence  $s(a)$  into a count vector  $\bar{x}_a = \{x_a^1, \dots, x_a^{|U|}\}$  where  $x_a^i$  are the counts of  
 10 the unique  $m$ -mers in  $U$ . For  $K$  clusters  $\bar{C} = \{C_1, \dots, C_K\}$ , the inertia is defined as

$$11 \quad arg \min_C \sum_{j=1}^K \sum_{\bar{x}_a \in C_j} |\bar{x}_a - \mu_j|^2$$

12 where  $\mu_j = \frac{1}{|C_j|} \sum_{\bar{x}_a \in C_j} \bar{x}_a$  is the mean of cluster  $j$ .

13 The recursive algorithm first partitions an MSA vertically into match and non-match intervals.  
 14 Match intervals are *collapsed* down to the single sequence they represent. Independently for  
 15 each non-match interval, the alignment slice is partitioned horizontally into clusters. The  
 16 same process is then applied to each induced sub-alignment until a maximum number of  
 17 recursion levels,  $r = 5$ , has been reached. For any remaining alignments, a node is added to  
 18 the local graph for each unique sequence. See Supplementary Animation 1 to see an  
 19 example of this algorithm. We name this algorithm Recursive Cluster and Collapse (RCC),  
 20 and implement in the `make_prg` repository (see Code Availability).

## 21 (w,k)-minimizers of graphs

22 We define (w,k)-minimizers of strings as in Li (2016)<sup>51</sup>. Let  $\phi : \Sigma^k \rightarrow \mathfrak{H}$  be a k-mer hash  
 23 function and let  $\pi : \Sigma^* \times \{0, 1\} \rightarrow \Sigma^*$  be defined such that  $\pi(s, 0) = s$  and  $\pi(s, 1) = \bar{s}$ , where  $\bar{s}$   
 24 is the reverse complement of  $s$ . Consider any integers  $k \geq w > 0$ . For window start position  
 25  $0 \leq j \leq |s| - w - k + 1$ , let

$$26 \quad T_j = \{\pi(s_{p,p+k}, r) : j \leq p < j + w, r \in \{0, 1\}\}$$

27 be the set of forward and reverse-complement k-mers of  $s$  in this window. We define a  
 28 (w,k)-minimizer to be any triple  $(h, p, r)$  such that

$$29 \quad h = \phi(\pi(s_{p,p+k}, r)) = \min\{\phi(t) : t \in T_j\}.$$

30 The set  $W(s)$  of (w,k)-minimizers for  $s$ , is the union of minimizers over such windows.

$$31 \quad W(s) = \bigcup_{0 \leq j \leq |s| - w - k + 1} \{(h, p, r) : h = \min\{\phi(t) : t \in T_j\}\}.$$

32 We extend this definition intuitively to an acyclic sequence graph  $G = (V, E)$ . Define  $|v|$  to be  
 33 the length of the sequence associated with node  $v \in V$  and let  $i = (v, a, b)$ ,  $0 \leq a \leq b \leq |v|$   
 34 represent the sequence interval  $[a, b]$  on  $v$ . We define a *path* in  $G$  by

$$35 \quad \bar{p} = \{(i_1, \dots, i_m) : (v_j, v_{j+1}) \in E \text{ and } b_j \equiv |v_j| \text{ for } 1 \leq j < m\}.$$

1 This matches the intuitive definition for a path in a sequence graph except that we allow the  
2 path to overlap only part of the sequence associated with the first and last nodes. We will  
3 use  $s_{\bar{p}}$  to refer to the sequence along the path  $\bar{p}$  in the graph.

4 Let  $\bar{p}$  be a path of length  $w+k-1$  in  $G$ . The string  $s_{\bar{p}}$  contains  $w$  consecutive  $k$ -mers for which  
5 we can find the  $(w,k)$ -minimizer(s) as before. We therefore define the  $(w,k)$ -minimizer(s) of  
6 the graph  $G$  to be the union of minimizers over all paths of length  $w+k-1$  in  $G$ :

$$7 \quad W(G) = \bigcup_{\bar{p} \in G: |\bar{p}|=w+k-1} \{(h, \bar{p}, r) : h = \min\{\phi(t) : t \in T_{\bar{p}}\}\}.$$

## 8 Local graph indexing with $(w,k)$ -minimizers

9 To find minimizers for a graph we use a streaming algorithm as described in Supplementary  
10 Algorithm 1. For each minimizer found, it simply finds the next minimizer(s) until the end of  
11 the graph has been reached.

12 Let  $walk(v, i, w, k)$  be a function which returns all vectors of  $w$  consecutive  $k$ -mers in  $G$   
13 starting at position  $i$  on node  $v$ . Suppose we have a vector of  $k$ -mers  $x$ . Let  $shift(x)$  be the  
14 function which returns all possible vectors of  $k$ -mers which extend  $x$  by one  $k$ -mer. It does  
15 this by considering possible ways to walk one letter in  $G$  from the end of the final  $k$ -mer of  $x$ .  
16 For a vector of  $k$ -mers of length  $w$ , the function  $minimize(x)$  returns the minimizing  $k$ -mers of  
17  $x$ .

18 We define  $K$  to be a  $k$ -mer graph with nodes corresponding to minimizers  $(h, \bar{p}, r)$ . We add  
19 edge  $(u,v)$  to  $K$  if there exists a path in  $G$  for which  $u$  and  $v$  are both minimizers and  $v$  is the  
20 first minimizer after  $u$  along the path. Let  $K \leftarrow add(s,t)$  denote the addition of nodes  $s$  and  $t$  to  
21  $K$  and the directed edge  $(s,t)$ . Let  $K \leftarrow add(s,T)$  denote the addition of nodes  $s$  and  $t \in T$  to  
22  $K$  as well as directed edges  $(s,t)$  for  $t \in T$ , and define  $K \leftarrow add(S,t)$  similarly.

23 The resulting PanRG index stores a map from each minimizing  $k$ -mer hash value to the  
24 positions in all local graphs where that  $(w,k)$ -minimizer occurred. In addition, we store the  
25 induced  $k$ -mer graph for each local graph.

## 26 Quasi-mapping reads

27 We infer the presence of PanRG loci in reads by quasi-mapping. For each read, a sketch of  
28  $(w,k)$ -minimizers is made, and these are queried in the index. For every  $(w,k)$ -minimizer  
29 shared between the read and a local graph in the PanRG index, we define a *hit* to be the  
30 coordinates of the minimizer in the read and local graph and whether it was found in the  
31 same or reverse orientation. We define clusters of hits from the same read, local graph and  
32 orientation if consecutive read coordinates are within a certain distance. If this cluster is of  
33 sufficient size, the locus is deemed to be present and we keep the hits for further analysis.  
34 Otherwise, they are discarded as noise. The default for this “sufficient size” is at least 10 hits  
35 and at least 1/5th the length of the shortest path through the  $k$ -mer graph (Nanopore) or the

1 number of k-mers in a read sketch (Illumina). Note that there is no requirement for all these  
2 hits to lie on a single path through the local graph. A further filtering step is therefore applied  
3 after the sequence at a locus is inferred to remove false positive loci, as indicated by low  
4 mean or median coverage along the inferred sequence by comparison with the global  
5 average coverage. This quasi-mapping procedure is described in pseudocode in  
6 Supplementary Algorithm 2.

## 7 Initial sequence approximation as a mosaic of references

8 For each locus identified as present in the set of reads, quasi-mapping provides (filtered)  
9 coverage information for nodes of the directed acyclic k-mer graph. We use these to  
10 approximate the sequence as a mosaic of references as follows. We model k-mer coverage  
11 with a negative binomial distribution and use the simplifying assumption that k-mers are read  
12 independently. Let  $\Theta$  be the set of possible paths through the k-mer graph, which could  
13 correspond to the true genomic sequence from which reads were generated. Let  $r + s$  be the  
14 number of times the underlying DNA was read by the machine, generating a k-mer coverage  
15 of  $s$ , and  $r$  instances where the k-mer was sequenced with errors. Let  $1 - p$  be the probability  
16 that a given k-mer was sequenced correctly. For any path  $\theta \in \Theta$ , let  $\{X_1, \dots, X_M\}$  be  
17 independent and identically distributed random variables with probability distribution  
18  $f(x_i, r, p) = \frac{\Gamma(r+s)}{\Gamma(r)!} p^r (1-p)^s$ , representing the k-mer coverages along this path. Since the mean  
19 and variance are  $\frac{(1-p)r}{p}$  and  $\frac{(1-p)r}{p^2}$  we solve for  $r$  and  $p$  using the observed k-mer coverage  
20 mean and variance across all k-mers in all graphs for the sample. Let  $D$  be the k-mer  
21 coverage data seen in the read dataset. We maximise the log-likelihood-inspired score  
22  $\hat{\theta} = \{arg \max_{\theta \in \Theta} l(\theta|D)\}$  where  $l(\theta|D) = \frac{1}{M} \sum_{i=1}^M \log f(s_i, r, p)$ , where  $s_i$  is the observed  
23 coverage of the  $i$ -th k-mer in  $\theta$ . By construction, the k-mer graph is directed and acyclic so  
24 this maximisation problem can be solved with a dynamic programming algorithm (for  
25 pseudocode, see Supplementary Algorithm 3).

26 For choices of  $w \leq k$  there is a unique sequence along the discovered path through the  
27 k-mer graph (except in rare cases within the first or last  $w-1$  bases). We use this closest  
28 mosaic of reference sequences as an initial approximation of the sample sequence.

## 29 *De novo* variant discovery

30 The first step in our implementation of local *de novo* variant discovery in genome graphs is  
31 finding the candidate regions of the graph that show evidence of dissimilarity from the  
32 sample's reads.

## 33 Finding candidate regions

1 The input required for finding candidate regions is a local graph,  $n$ , within the PanRG, the  
2 maximum likelihood path of both sequence and k-mers in  $n$ ,  $lmp_n$  and  $kmp_n$  respectively, and  
3 a padding size  $w$  for the number of positions surrounding the candidate region to retrieve.

4 We define a candidate region,  $r$ , as an interval within  $n$  where coverage on  $lmp_n$  is less than  
5 a given threshold,  $c$ , for more than  $l$  and less than  $m$  consecutive positions.  $m$  acts to restrict  
6 the size of variants we are able to detect. If set too large, the following steps become much  
7 slower due to the combinatorial expansion of possible paths.

8 For a given read,  $s$ , that has a mapping to  $r$  we define  $s_r$  to be the subsequence of  $s$  that  
9 maps to  $r$ , including an extra  $w$  positions either side of the mapping. We define the pileup  $P_r$   
10 as the set of all  $s_r \in r$ .

## 11 Enumerating paths through candidate regions

12 For  $r \in R$ , where  $R$  is the set of all candidate regions, we construct a de Bruijn graph  $G_r$   
13 from the pileup  $P_r$  using the GATB library<sup>52</sup>.  $A_L$  and  $A_R$  are defined as sets of k-mers to the  
14 left and right of  $r$  in the local graph. They are anchors to allow re-insertion of new sequences  
15 found by *de novo* discovery into the local graph. If we cannot find an anchor on both sides,  
16 then we abandon *de novo* discovery for  $r$ . We use sets of k-mers for  $A_L$  and  $A_R$ , rather than  
17 a single anchor k-mer, to provide redundancy in the case where sequencing errors cause  
18 the absence of some k-mers in  $G_r$ . Once  $G_r$  is built, we define the start anchor k-mer,  $a_L$ ,  
19 as the first  $a_L \in A_L \wedge a_L \in G_r$ . Likewise, we define the end anchor k-mer,  $a_R$ , as  
20 the first  $a_R \in A_R \wedge a_R \in G_r$ .

21  $T_r$  is the spanning tree obtained by performing depth-first search (DFS) on  $G_r$ , beginning  
22 from node  $a_L$ . We define  $p_r$  as a path, from the root node  $a_L$  of  $T_r$  and ending at node  $a_R$ ,  
23 which fulfils the two conditions: 1)  $p_r$  is shorter than the maximum allowed path length. 2)  
24 No more than  $k$  nodes along  $p_r$  have coverage  $< f \times e_r$ , where  $e_r$  is the expected k-mer  
25 coverage for  $r$  and  $f$  is  $n_r \times s$ , where  $n_r$  is the number of iterations of path enumeration for  $r$   
26 and  $s$  is a step size (0.1 by default).

27  $V_r$  is the set of all  $p_r$ . If  $|V_r|$  is greater than a predefined threshold, then we have too many  
28 candidate paths, and we decide to filter more aggressively:  $f$  is incremented by  $s$  - effectively  
29 requiring more coverage for each  $p_r$  - and  $V_r$  is repopulated. If  $f > 1.0$  then *de novo*  
30 discovery is abandoned for  $r$ .

## 31 Pruning the path-space in a candidate region

32 As we operate on both accurate and error-prone sequencing reads, the number of valid  
33 paths in  $G_r$  can be very large. Primarily, this is due to cycles that can occur in  $G_r$  and  
34 exploring paths that will never reach our required end anchor  $a_R$ . In order to reduce the

1 path-space within  $G_r$  we prune paths based on multiple criteria. Critically, this pruning  
2 happens at each step of the graph walk (path-building).  
3 We used a distance-based optimisation based on Rizzi et al <sup>53</sup>. In addition to  $T_r$ , obtained  
4 by performing DFS on  $G_r$ , we produce a distance map  $D_r$  that results from running  
5 reversed breadth-first search (BFS) on  $G_r$ , beginning from node  $a_R$ . We say reversed BFS  
6 as we explore the predecessors of each node, rather than the successors.  $D_r$  is  
7 implemented as a binary search tree where each node in the tree represents a k-mer in  $G_r$   
8 that is reachable from  $a_R$  via reversed BFS. Each node additionally has an integer attached  
9 to it that describes the distance from that node to  $a_R$ .  
10 We can use  $D_r$  to prune the path-space by 1) for each node  $n \in p_r$ , we require  $n \in D_r$  and  
11 2) requiring  $a_R$  be reached from  $n$  in, at most,  $i$  nodes, where  $i$  is defined as the maximum  
12 allowed path length minus the number of nodes walked to reach  $n$ .  
13 If one of these conditions is not met, we abandon  $p_r$ . The advantage of this pruning process  
14 is that we never explore paths that will not reach our required endpoint within the maximum  
15 allowed path length and when caught in a cycle, we abandon the path once we have made  
16 too many iterations around the cycle.

## 17 Graph-based genotyping and optimal reference construction for 18 multi-genome comparison

19 We use graph-based genotyping to output a comparison of samples in a VCF. A path  
20 through the graph is selected to be the reference sequence, and graph variation is described  
21 with respect to this reference. The chromosome field then details the local graph and the  
22 position field gives the position within the chosen reference sequence for possible variant  
23 alleles. The reference path for each local graph is chosen to be maximally close to the set of  
24 sample mosaic paths. This is achieved by reusing the mosaic path finding algorithm detailed  
25 in Supplementary Algorithm 3 on a copy of the k-mer graph with coverages incremented  
26 along each sample mosaic path, and a modified probability function defined such that the  
27 probability of a node is proportional to the number of samples covering it. This results in an  
28 optimal path, which is used as the VCF reference for the multi-sample VCF file.

29 For each sample and site in the VCF file, the mean forward and reverse coverage on k-mers  
30 tiling alleles is calculated. A likelihood is then independently calculated for each allele based  
31 on a Poisson model. An allele  $A$  in a site is called if: 1)  $A$  is on the sample mosaic path (i.e.  
32 it is on the maximum likelihood path for that sample); 2)  $A$  is the most likely allele to be  
33 called based on the previous Poisson model. Every allele not in the sample mosaic path will  
34 not satisfy 1) and will thus not be called. In the uncommon event where an allele satisfies 1),  
35 but not 2), we have an incompatibility between the global and the local choices, and then the  
36 site is genotyped as null.

## 1 Comparison of variant-callers on a diverse set of *E. coli*

### 2 Sample selection

3 We used a set of 20 diverse *E. coli* samples for which matched Nanopore and Illumina data  
4 and a high-quality assembly were available. These are distributed across 4 major  
5 phylogroups of *E. coli* as shown in Figure 4. Of these, 16 were isolated from clinical  
6 infections and rectal screening swabs in ICU patients in an Australian hospital<sup>54</sup>. One is the  
7 reference strain CFT073 that was resequenced and assembled by the REHAB consortium<sup>55</sup>.  
8 One is from an ST216 cardiac ward outbreak (identifier: H131800734); the Illumina data was  
9 previously obtained<sup>56</sup> and we did the Nanopore sequencing (see below). The two final  
10 samples were obtained from Public Health England: one is a Shiga-toxin encoding *E. coli*  
11 (we used the identifier O63)<sup>57</sup>, and the other an enteroaggregative *E. coli* (we used the  
12 identifier ST38)<sup>58</sup>. Coverage data for these samples can be found in Supplementary Table 1.

### 13 PanRG construction

14 MSAs for gene clusters curated with PanX<sup>27</sup> from 350 RefSeq assemblies were downloaded  
15 from <http://pangenome.de> on 3rd May 2018. MSAs for intergenic region clusters based on  
16 228 *E. coli* ST131 genome sequences were previously generated with Piggy<sup>28</sup> for their  
17 publication. The PanRG was built using *make\_prg*. Two loci (GC00000027\_2 and  
18 GC00004221) out of 37,428 were excluded because the combination of Clustal Omega and  
19 *make\_prg* did not complete in reasonable time (~24 hours) once *de novo* variants were  
20 added.

### 21 Nanopore sequencing of sample H131800734

22 DNA was extracted using a Blood & Cell Culture DNA Midi Kit (Qiagen, Germany) and  
23 prepared for Nanopore sequencing using kits EXP-NBD103 and SQK-LSK108. Sequencing  
24 was performed on a MinION Mk1 Shield device using a FLO-MIN106 R9.4 Spoton flowcell  
25 and MinKNOW version 1.7.3, for 48 hours.

### 26 Nanopore basecalling

27 Recent improvements to the accuracy of Nanopore reads have been largely driven by  
28 improvements in basecalling algorithms<sup>59</sup>. All Nanopore data was basecalled with the  
29 methylation-aware, high-accuracy model provided with the proprietary guppy basecaller  
30 (version 3.4.5). In addition, 4 samples were basecalled with the default (methylation  
31 unaware) model for comparison (see Figure 5). Demultiplexing of the subsequent basecalled  
32 data was performed using the same version of the guppy software suite with barcode kits  
33 EXP-NBD104 and EXP-NBD114 and an option to trim the barcodes from the output.

## 1 Phylogenetic tree construction

2 Chromosomes were aligned using *MAFFT*<sup>60</sup> v7.467 as implemented in *Parsnp*<sup>61</sup> v1.5.3.  
3 *Gubbins* v2.4.1 was used to filter for recombination (default settings) and phylogenetic  
4 construction was carried out using *RAxML*<sup>62</sup> v8.2.12 (GTR + GAMMA substitution model, as  
5 implemented in *Gubbins*<sup>63</sup>).

## 6 Reference selection for mapping-based callers

7 A set of references was chosen for testing single-reference variant callers using two  
8 standard approaches, as follows. First, a phylogeny was built containing our 20 samples and  
9 243 reference genomes from RefSeq. Then, for each of our 20 samples, the nearest RefSeq  
10 *E. coli* reference was found using *Mash*<sup>64</sup>. Second, for each of the 20 samples, the nearest  
11 RefSeq reference in the phylogeny was manually selected; sometimes one RefSeq  
12 assembly was the closest to more than one of the 20. At an earlier stage of the project there  
13 had been another sample (making a total of 21) in phylogroup B1; this was discarded when  
14 it failed quality filters (data not shown). Despite this, the *Mash*/manual selected reference  
15 genomes were left in the set of mapping references, to evaluate the impact of mapping to a  
16 reference in a different phylogroup to all 20 of our samples.

## 17 Construction of truth assemblies

18 16/20 samples were obtained with matched Illumina and Nanopore data and a hybrid  
19 assembly. Sample H131800734 was assembled using the hybrid assembler *Unicycler*<sup>65</sup> with  
20 PacBio and Illumina reads followed by polishing with the PacBio reads using *Racon*<sup>66</sup>, and  
21 finally with Illumina reads using *Pilon*<sup>67</sup>. A small 1kb artifactual contig was removed from the  
22 H131800734 assembly due to low quality and coverage.

23 In all cases we mapped the Illumina data to the assembly, and masked all positions where  
24 the pileup of Illumina reads did not support the assembly.

## 25 Construction of a comprehensive and filtered truth set of pairwise SNPs

26 All pairwise comparisons of the 20 truth assemblies were performed with *varifier*  
27 (<https://github.com/iqbal-lab-org/varifier>), using subcommand *make\_truth\_vcf*. In summary,  
28 *varifier* compares two given genomes (referenced as G1 and G2) twice - first using *dnadiff*<sup>68</sup>  
29 and then using *minimap2/paftools*<sup>51</sup>. The two output sets of pairwise SNPs are then joined  
30 and filtered. We create one sequence probe for each allele (a sequence composed of the  
31 allele and 50 bases of flank on either side taken from G1) and then map both to G2 using  
32 *minimap2*. We then evaluate these mappings to verify if the variant found is indeed correct  
33 (TP) or not (FP) as follows. If the mapping quality is zero, the variant is discarded to avoid  
34 paralogs/duplicates/repeats that are inherently hard to assess. We then check for  
35 mismatches in the allele after mapping and confirm that the called allele is the better match.

## 36 Constructing a set of ground truth pan-genome variants

37 When seeking to construct a truth set of all variants within a set of bacterial genomes, there  
38 is no universal coordinate system. We start by taking all pairs of genomes and finding the

1 variants between them, and then need to deduplicate them - e.g. when a variant between  
2 genomes 1 and 2 is the same as a variant between genomes 3 and 4, they should be  
3 identified; we define “the same” in terms of genome, coordinate and allele. An allele  $A$  in a  
4 position  $P_A$  of a chromosome  $C_A$  in a genome  $G_A$  is defined as a triple  $A = (G_A, C_A, P_A)$ .  
5 A pairwise variant  $PwV = \{A_1, A_2\}$  is defined as a pair of alleles that describes a variant  
6 between two genomes, and a pan-genome variant  $PgV = \{A_1, A_2, \dots, A_n\}$  is defined as a set  
7 of two or more alleles that describes the same variant between two or more genomes. A  
8 pan-genome variant  $PgV$  can also be defined as a set of pairwise variants  
9  $PgV = \{PwV_1, PwV_2, \dots, PwV_n\}$ , as we can infer the set of alleles of  $PgV$  from the pairs  
10 of alleles in all these pairwise variants. Note that pan-genome variants are thus able to  
11 represent rare and core variants. Given a set of pairwise variants, we seek a set of  
12 pan-genome variants satisfying the following properties:  
13 1. [Surjection]:  
14 a. each pairwise variant is in exactly one pan-genome variant;  
15 b. a pan-genome variant contains at least one pairwise variant;  
16 2. [Transitivity]: if two pairwise variants  $PwV_1$  and  $PwV_2$  share an allele, then  $PwV_1$   
17 and  $PwV_2$  are in the same pan-genome variant  $PgV$ ;

18 We model the above problem as a graph problem. We represent each pairwise variant as a  
19 node in an undirected graph  $G$ . There is an edge between two nodes  $n_1$  and  $n_2$  if  $n_1$  and  
20  $n_2$  share an allele. Each component (maximal connected subgraph) of  $G$  then defines a  
21 pan-genome variant, built from the set of pairwise variants in the component, satisfying all  
22 the properties previously described. Therefore, the set of components of  $G$  defines the set  
23 of pan-genome variants  $P$ . However, a pan-genome variant in  $P$  could: i) have more than  
24 one allele stemming from a single genome, due to a duplication/repeat; ii) represent biallelic  
25 , triallelic or tetraallelic SNPs/indels. For this evaluation, we chose to have a smaller, but more  
26 reliable set of pan-genome variants, and thus we filtered  $P$  by restricting it to the set of  
27 pan-genome variants  $P'$  defined by the variants  $PgV \in P$  such that: i)  $PgV$  has at most  
28 one allele stemming from each genome; ii)  $PgV$  is a biallelic SNP.  $P'$  is the set of 618,305  
29 ground truth filtered pan-genome variants that we extracted by comparing and deduplicating  
30 the pairwise variants present in our 20 samples, and that we use to evaluate the recall of all  
31 the tools in this paper. Supplementary Figure 11 shows an example summarising the  
32 described process of building pan-genome variants from a set of pairwise variants.

### 33 Subsampling read data and running all tools

34 All read data was randomly subsampled to 100x coverage using *rasusa* - the pipeline is  
35 available at <https://github.com/iqbal-lab-org/subsampler>. A *snakemake*<sup>69</sup> pipeline to run the  
36 *pandora* workflow with and without *de novo* discovery (see Figure 2d) is available at  
37 [https://github.com/iqbal-lab-org/pandora\\_workflow](https://github.com/iqbal-lab-org/pandora_workflow). A *snakemake* pipeline to run *snippy*,  
38 *SAMtools*, *nanopolish* and *medaka* on all pairwise combinations of 20 samples and 24  
39 references is available at [https://github.com/iqbal-lab-org/variant\\_callers\\_pipeline](https://github.com/iqbal-lab-org/variant_callers_pipeline).

## 1 Evaluating VCF files

### 2 Calculating precision

3 Given a variant/VCF call made by any of the evaluated tools, where the input were reads  
4 from a sample (or several samples, in the case of *pandora*) and a reference sequence (or a  
5 PanRG, in the case of *pandora*), we perform the following steps to assess how correct a call  
6 is:

- 7 1. Construct a probe for the called allele, consisting of the sequence of the allele  
8 flanked by 150bp on both sides from the reference sequence. This reference  
9 sequence is one of the 24 chosen references for *snippy*, *SAMtools*, *nanopolish* and  
10 *medaka*; or the multi-sample inferred VCF reference for *pandora*;
- 11 2. Map the probe to the sample sequence using *BWA-MEM*<sup>70</sup>;
- 12 3. Remove multi-mappings by looking at the Mapping Quality (MAPQ) measure<sup>30</sup> of the  
13 SAM records. If the probe is mapped uniquely, then its mapping passes the filter. If  
14 there are multiple mappings for the probe, we select the mapping  $m_1$  with the highest  
15 MAPQ if the difference between its MAPQ and the second highest MAPQ exceeds  
16 10. If  $m_1$  does not exist, then there are at least two mappings with the same MAPQ,  
17 and it is ambiguous to choose which one to evaluate. In this case, we prefer to be  
18 conservative and filter this call (and all its related mappings) out of the evaluation;
- 19 4. We further remove calls mapping to masked regions of the sample sequence, in  
20 order to not evaluate calls lying on potentially misassembled regions;
- 21 5. Now we evaluate the mapping, giving the call a continuous precision score between  
22 0 and 1. If the mapping does not cover the whole called allele, we give a score of 0.  
23 Otherwise, we look only at the alignment of the called allele (i.e. we ignore the  
24 flanking sequences alignment), and give a score of: number of matches / alignment  
25 length.

26 Finally, we compute the precision for the tool by summing the score of all evaluated calls and  
27 dividing by the number of evaluated calls. Note that here we evaluate all types of variants,  
28 including SNPs and indels.

### 29 Calculating recall

30 We perform the following steps to calculate the recall of a tool:

- 31 1. Apply the VCF calls to the associated reference using the VCF consensus builder  
32 ([https://github.com/leoisl/vcf\\_consensus\\_builder](https://github.com/leoisl/vcf_consensus_builder)), creating a mutated reference with  
33 the variants identified by the tool;
- 34 2. Build probes for each allele of each pan-genome variant previously computed (see  
35 Section “Constructing a set of ground truth pan-genome variants”);
- 36 3. Map all pan-genome variants’ probes to the mutated reference using *BWA-MEM*;
- 37 4. Evaluate each probe mapping, which is classified as a TP only if all bases of the  
38 allele were correctly mapped to the mutated reference. In the uncommon case where  
39 a probe multimaps, it is enough that one of the mappings are classified as TP;

1 5. Finally, as we now know for each pan-genome variant which of its alleles were found,  
2 we calculate both the pan-variant recall and the average allelic recall as per Section  
3 “*Pandora detects rare variation inaccessible to single-reference methods*”.

#### 4 Filters

5 Given a VCF file with likelihoods for each genotype, the genotype confidence is defined as  
6 the log likelihood of the maximum likelihood genotype, minus the log likelihood of the next  
7 best genotype. Thus a confidence of zero means all alleles are equally likely, and high  
8 quality calls have higher confidences. In the recall/error rate plots of Figure 5 and Figures  
9 6a,b, each point corresponds to the error rate and recall computed as previously described,  
10 on a genotype confidence (gt-conf) filtered VCF file with a specific threshold for minimum  
11 confidence.

12 We also show the same plot with further filters applied in Supplementary Figure 1. The filters  
13 were as follows. For Illumina data: for *pandora*, a minimum coverage filter of 5x, a strand  
14 bias filter of 0.05 (minimum 5% of reads on each strand), and a gaps filter of 0.8 were  
15 applied. The gaps filter means at least 20% the minimizer k-mers on the called allele must  
16 have coverage above 10% of the expected depth. As *snippy* has its own internal filtering, no  
17 filters were applied. For *SAMtools*, a minimum coverage filter of 5x was used. For Nanopore  
18 data: for *pandora*, a minimum coverage filter of 10x, a strand bias filter of 0.05, and a gaps  
19 filter of 0.6 were used. For *nanopolish*, we applied a coverage filter of 10x. We were unable  
20 to apply a minimum coverage filter to a *medaka* due to a software bug that prevents  
21 annotating the VCF file with coverage information.

#### 22 Locus presence and distance evaluation

23 For all loci detected as present in at least one sample by *pandora*, we mapped the  
24 multi-sample inferred reference to all 20 sample assemblies and 24 references, to identify  
25 their true locations. To be confident of these locations, we employed a strict mapping using  
26 *bowtie2*<sup>71</sup> and requiring end-to-end alignments. From the mapping of all loci to all samples,  
27 we computed a truth locus presence-absence matrix, and compared it with *pandora*'s locus  
28 presence-absence matrix, classifying each *pandora* locus call as true/false positive/negative.  
29 Supplementary Figure 3 shows these classifications split by locus length. Having the location  
30 of all loci in all the 20 sample assemblies and the 24 references, we then computed the edit  
31 distance between them.

#### 32 Reproducibility

33 All input data for our analyses, including PanX's and Piggy's MSAs, PanRG, reference  
34 sequences, and sample data are publicly available (see Section “*Data availability*”).  
35 *Pandora*'s code, as well as all code needed to reproduce this analysis are also publicly  
36 available (see Section “*Code availability*”). Software environment reproducibility is achieved  
37 using Python virtual environments if all dependencies and source code are in Python, and  
38 using Docker<sup>72</sup> containers run with Singularity<sup>73</sup> otherwise. The exact commit/version of all  
39 repositories used to obtain the results in this paper can be retrieved with the git branch or  
40 tag *pandora\_paper\_tag1*.

## 1 Data availability

- 2 ● Gene MSAs from PanX, and intergenic MSAs from Piggy:  
3 [doi.org/10.6084/m9.figshare.13204163](https://doi.org/10.6084/m9.figshare.13204163);
- 4 ● *E. Coli* PanRG: [doi.org/10.6084/m9.figshare.13204172](https://doi.org/10.6084/m9.figshare.13204172);
- 5 ● Accession identifiers or Figshare links for the sample and reference assemblies, and  
6 Illumina and Nanopore reads are listed in Section D of the Supplementary file;
- 7 ● Input packages containing all data to reproduce both the 4- and 20-way analyses  
8 described in the Results section are also available in Section D of the Supplementary  
9 file.

## 10 Code availability

- 11 ● *make\_prg* (RCC graph construction algorithm): [https://github.com/rmcolq/make\\_prg](https://github.com/rmcolq/make_prg)
- 12 ● *pandora*: <https://github.com/rmcolq/pandora>
- 13 ● *varifier*: <https://github.com/iqbal-lab-org/varifier>
- 14 ● Pangenome variations pipeline taking a set of assemblies and returning a set of  
15 filtered pan-genome variants: [https://github.com/iqbal-lab-org/pangenome\\_variations](https://github.com/iqbal-lab-org/pangenome_variations)
- 16 ● *pandora* workflow: [https://github.com/iqbal-lab-org/pandora\\_workflow](https://github.com/iqbal-lab-org/pandora_workflow)
- 17 ● Run *snippy*, *samtools*, *nanopolish* and *medaka* pipeline:  
18 [https://github.com/iqbal-lab-org/variant\\_callers\\_pipeline](https://github.com/iqbal-lab-org/variant_callers_pipeline)
- 19 ● 4- and 20-way evaluation pipeline (recall/error rate curves etc):  
20 [https://github.com/iqbal-lab-org/pandora\\_paper\\_roc](https://github.com/iqbal-lab-org/pandora_paper_roc)
- 21 ● Locus presence and distance from reference pipeline:  
22 [https://github.com/iqbal-lab-org/pandora\\_gene\\_distance](https://github.com/iqbal-lab-org/pandora_gene_distance)
- 23 ● A master repository to reproduce everything in this paper, marshalling all of the  
24 above: [https://github.com/iqbal-lab-org/paper\\_pandora2020\\_analyses](https://github.com/iqbal-lab-org/paper_pandora2020_analyses)

25 Although all containers are hosted on <https://hub.docker.com/> (for details, see  
26 [https://github.com/iqbal-lab-org/paper\\_pandora2020\\_analyses/blob/master/scripts/pull\\_containers/pull\\_containers.sh](https://github.com/iqbal-lab-org/paper_pandora2020_analyses/blob/master/scripts/pull_containers/pull_containers.sh)), and are downloaded automatically during the pipelines' execution,  
27 we also provide Singularity<sup>73</sup> containers (converted from Docker containers) at  
28 [doi.org/10.6084/m9.figshare.13204169](https://doi.org/10.6084/m9.figshare.13204169).

30 Frozen packages with all the code repositories for *pandora* and the analysis framework can  
31 be found at [doi.org/10.6084/m9.figshare.13204214](https://doi.org/10.6084/m9.figshare.13204214).

## 32 Author contributions

33 RMC designed and implemented the fundamental data structures, and RCC, map and  
34 compare algorithms. MBH designed and implemented the *de novo* variant discovery  
35 component. LL optimised the codebase. RMC, MBH, LL designed and implemented (several  
36 iterations of) the evaluation pipeline, one component of which was written by MH. BL  
37 reimplemented and improved the RCC codebase. JH,SG,LP sequenced 18/20 of the

1 samples. LWR, MBH, LL, KM, ZI analysed and visualised the 20-way data. ZI designed the  
2 study. ZI and RMC wrote the bulk of the paper, LL and MBH wrote sections, and all authors  
3 read and improved drafts.

#### 4 Competing interest declaration

5 The authors declare no competing interests.

#### 6 Acknowledgements

7 We are grateful to the REHAB consortium (<https://modmedmicro.nsms.ox.ac.uk/rehab/>) and  
8 the Transmission of Carbapenemase-producing Enterobacteriaceae (TRACE) study  
9 investigators for sharing sequencing data (for CFT073 and H131800734) in support of this  
10 work. We would like to thank Sion Bayliss and Ed Thorpe for discussions and help with  
11 Piggy. We are grateful to Kelly Wyres for sharing sequence data for the Australian samples,  
12 and to Tim Dallman and David Greig for sharing their data from Public Health England. We  
13 would like to thank the following for helpful conversations during the prolonged genesis of  
14 this project: Gil McVean, Derrick Crook, Eduardo Rocha, Bill Hanage, Ed Feil, Sion Bayliss,  
15 Ed Thorpe, Richard Neher, Camille Marchet, Rayan Chikhi, Kat Holt, Claire Gorrie, Rob  
16 Patro, Fatemeh Almodaresi, Nicole Stoesser, Liam Shaw, Phelim Bradley, and Sorina  
17 Maciuca. RMC was funded by a Wellcome Trust PhD studentship (105279/Z/14/Z), and ZI  
18 was partially funded by a Wellcome Trust/Royal Society Sir Henry Dale Fellowship  
19 (102541/Z/13/Z).

#### 20 References

- 21 1. Lynch, M. *et al.* Genetic drift, selection and the evolution of the mutation rate. *Nat. Rev.*  
22 *Genet.* **17**, 704–714 (2016).
- 23 2. Didelot, X. & Maiden, M. C. J. Impact of recombination on bacterial evolution. *Trends*  
24 *Microbiol.* **18**, 315–322 (2010).
- 25 3. Rocha, E. P. C. Neutral Theory, Microbial Practice: Challenges in Bacterial Population  
26 Genetics. *Mol. Biol. Evol.* **35**, 1338–1347 (2018).
- 27 4. Fraser, C., Alm, E. J., Polz, M. F., Spratt, B. G. & Hanage, W. P. The Bacterial Species  
28 Challenge: Making Sense of Genetic and Ecological Diversity. *Science* **323**, 741–746  
29 (2009).
- 30 5. Mira, A., Ochman, H. & Moran, N. A. Deletional bias and the evolution of bacterial

- 1 genomes. *Trends Genet.* **17**, 589–596 (2001).
- 2 6. Domingo-Sananes, M. R. & McInerney, J. Selection-based model of prokaryote  
3 pangenomes | bioRxiv. <https://www.biorxiv.org/content/10.1101/782573v1>.
- 4 7. Gordienko, E. N., Kazanov, M. D. & Gelfand, M. S. Evolution of Pan-Genomes of  
5 *Escherichia coli*, *Shigella* spp., and *Salmonella enterica*. *J. Bacteriol.* **195**, 2786–2792  
6 (2013).
- 7 8. Lobkovski, A., Wolf, Y. & Koonin, Eugene. Gene Frequency Distributions Reject a  
8 Neutral Model of Genome Evolution | Genome Biology and Evolution | Oxford Academic.  
9 <https://academic.oup.com/gbe/article/5/1/233/732669>.
- 10 9. Bolotin, E. & Hershberg, R. Gene Loss Dominates As a Source of Genetic Variation  
11 within Clonal Pathogenic Bacterial Species. *Genome Biol. Evol.* **7**, 2173–2187 (2015).
- 12 10. Haegeman, B. & Weitz, J. S. A neutral theory of genome evolution and the frequency  
13 distribution of genes. *BMC Genomics* **13**, 196 (2012).
- 14 11. Hadfield, J. *et al.* Phandango: an interactive viewer for bacterial population genomics.  
15 *Bioinformatics* **34**, 292–293 (2018).
- 16 12. Garrison, E. *et al.* Variation graph toolkit improves read mapping by representing genetic  
17 variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
- 18 13. Maciuca, S., Elias, C. del O., McVean, G. & Iqbal, Z. A natural encoding of genetic  
19 variation in a Burrows-Wheeler Transform to enable mapping and genome inference.  
20 *bioRxiv* 059170 (2016) doi:10.1101/059170.
- 21 14. Eggertsson, H. P. *et al.* GraphTyper enables population-scale genotyping using  
22 pangenome graphs. *Nat. Genet.* **49**, 1654–1660 (2017).
- 23 15. Eggertsson, H. P. *et al.* GraphTyper2 enables population-scale genotyping of structural  
24 variation using pangenome graphs. *Nat. Commun.* **10**, 5402 (2019).
- 25 16. Rautiainen, M. & Marschall, T. GraphAligner: Rapid and Versatile Sequence-to-Graph  
26 Alignment. *bioRxiv* 810812 (2019) doi:10.1101/810812.

- 1 17. Schneeberger, K. *et al.* Simultaneous alignment of short reads against multiple  
2 genomes. *Genome Biol.* **10**, R98 (2009).
- 3 18. Rabbani, L., Müller, J. & Weigel, D. An Algorithm to Build a Multi-genome Reference.  
4 *bioRxiv* 2020.04.11.036871 (2020) doi:10.1101/2020.04.11.036871.
- 5 19. The Computational Pan-Genomics Consortium. Computational pan-genomics: status,  
6 promises and challenges | Briefings in Bioinformatics | Oxford Academic.  
7 <https://academic.oup.com/bib/article/19/1/118/2566735>.
- 8 20. Rautiainen, M. & Marschall, T. Aligning sequences to general graphs in  $O(V + mE)$  time.  
9 *bioRxiv* 216127 (2017) doi:10.1101/216127.
- 10 21. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158  
11 (2011).
- 12 22. Roberts, M., Hayes, W., Hunt, B. R., Mount, S. M. & Yorke, J. A. Reducing storage  
13 requirements for biological sequence comparison. *Bioinformatics* **20**, 3363–3369 (2004).
- 14 23. Touchon, M. *et al.* Phylogenetic background and habitat drive the genetic diversification  
15 of *Escherichia coli*. *PLOS Genet.* **16**, e1008866 (2020).
- 16 24. Touchon, M. *et al.* Organised Genome Dynamics in the *Escherichia coli* Species Results  
17 in Highly Diverse Adaptive Paths. *PLOS Genet.* **5**, e1000344 (2009).
- 18 25. Decano, A. G. & Downing, T. An *Escherichia coli* ST131 pangenome atlas reveals  
19 population structure and evolution across 4,071 isolates. *Sci. Rep.* **9**, 17394 (2019).
- 20 26. Rasko, D. A. *et al.* The Pangenome Structure of *Escherichia coli*: Comparative Genomic  
21 Analysis of *E. coli* Commensal and Pathogenic Isolates. *J. Bacteriol.* **190**, 6881–6893  
22 (2008).
- 23 27. Ding, W., Baumdicker, F. & Neher, R. A. panX: pan-genome analysis and exploration.  
24 *Nucleic Acids Res.* **46**, e5–e5 (2018).
- 25 28. Thorpe, H. A., Bayliss, S. C., Sheppard, S. K. & Feil, E. J. Piggy: a rapid, large-scale  
26 pan-genome analysis tool for intergenic regions in bacteria. *GigaScience* **7**, (2018).

- 1 29. Clermont, O., Christenson, J. K., Denamur, E. & Gordon, D. M. The Clermont  
2 Escherichia coli phylo-typing method revisited: improvement of specificity and detection  
3 of new phylo-groups. *Environ. Microbiol. Rep.* **5**, 58–65 (2013).
- 4 30. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
5 2078–2079 (2009).
- 6 31. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing.  
7 *ArXiv12073907 Q-Bio* (2012).
- 8 32. Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de  
9 novo using only nanopore sequencing data. *Nat. Methods* **12**, 733–735 (2015).
- 10 33. Sievers, F. & Higgins, D. G. Clustal Omega for making accurate alignments of many  
11 protein sequences. *Protein Sci. Publ. Protein Soc.* **27**, 135–145 (2018).
- 12 34. Louca, S., Mazel, F., Doebeli, M. & Parfrey, L. W. A census-based estimate of Earth's  
13 bacterial and archaeal diversity. *PLOS Biol.* **17**, e3000106 (2019).
- 14 35. Brockhurst, M. A. *et al.* The Ecology and Evolution of Pangenomes. *Curr. Biol. CB* **29**,  
15 R1094–R1103 (2019).
- 16 36. Harrison, E. & Brockhurst, M. A. Plasmid-mediated horizontal gene transfer is a  
17 coevolutionary process. *Trends Microbiol.* **20**, 262–267 (2012).
- 18 37. Harrison, E. *et al.* Rapid compensatory evolution promotes the survival of conjugative  
19 plasmids. *Mob. Genet. Elem.* **6**, e1179074 (2016).
- 20 38. Loftie-Eaton, W. *et al.* Compensatory mutations improve general permissiveness to  
21 antibiotic resistance plasmids. *Nat. Ecol. Evol.* **1**, 1354–1363 (2017).
- 22 39. Gori, A. *et al.* Pan-GWAS of Streptococcus agalactiae Highlights Lineage-Specific  
23 Genes Associated with Virulence and Niche Adaptation. *mBio* **11**, (2020).
- 24 40. Bonnet, R. Growing Group of Extended-Spectrum  $\beta$ -Lactamases: the CTX-M Enzymes.  
25 *Antimicrob. Agents Chemother.* **48**, 1–14 (2004).
- 26 41. Louwen, R., Staals, R. H. J., Endtz, H. P., Baarlen, P. van & Oost, J. van der. The Role

- 1 of CRISPR-Cas Systems in Virulence of Pathogenic Bacteria. *Microbiol. Mol. Biol. Rev.*  
2 **78**, 74–88 (2014).
- 3 42. Horvath, P. *et al.* Diversity, Activity, and Evolution of CRISPR Loci in *Streptococcus*  
4 *thermophilus*. *J. Bacteriol.* **190**, 1401–1412 (2008).
- 5 43. Pritt, J., Chen, N.-C. & Langmead, B. FORGe: prioritizing variants for graph genomes.  
6 *Genome Biol.* **19**, 220 (2018).
- 7 44. Norri, T., Cazaux, B., Kosolobov, D. & Mäkinen, V. Linear time minimum segmentation  
8 enables scalable founder reconstruction. *Algorithms Mol. Biol.* **14**, 12 (2019).
- 9 45. Horesh, G. *et al.* A comprehensive and high-quality collection of *E. coli* genomes and  
10 their genes. *bioRxiv* 2020.09.21.293175 (2020) doi:10.1101/2020.09.21.293175.
- 11 46. Lees, J. A. *et al.* Sequence element enrichment analysis to determine the genetic basis  
12 of bacterial phenotypes. *Nat. Commun.* **7**, 1–8 (2016).
- 13 47. Earle, S. G. *et al.* Identifying lineage effects when controlling for population structure  
14 improves power in bacterial association studies. *Nat. Microbiol.* **1**, 1–8 (2016).
- 15 48. Jaillard, M. *et al.* A fast and agnostic method for bacterial genome-wide association  
16 studies: Bridging the gap between k-mers and genetic events. *PLOS Genet.* **14**,  
17 e1007758 (2018).
- 18 49. Sheppard, S. K., Jolley, K. A. & Maiden, M. C. J. A Gene-By-Gene Approach to Bacterial  
19 Population Genomics: Whole Genome MLST of *Campylobacter*. *Genes* **3**, 261–277  
20 (2012).
- 21 50. MacQueen, J. Some methods for classification and analysis of multivariate observations.  
22 in (The Regents of the University of California, 1967).
- 23 51. Li, H. Minimap and miniasm: fast mapping and de novo assembly for noisy long  
24 sequences. *Bioinformatics* **32**, 2103–2110 (2016).
- 25 52. Drezen, E. *et al.* GATB: Genome Assembly & Analysis Tool Box. *Bioinformatics* **30**,  
26 2959–2961 (2014).

- 1 53. Rizzi, R., Sacomoto, G. & Sagot, M.-F. Efficiently Listing Bounded Length st-Paths. in  
2 *Combinatorial Algorithms* (eds. Jan, K., Miller, M. & Froncek, D.) 318–329 (Springer  
3 International Publishing, 2015). doi:10.1007/978-3-319-19315-1\_28.
- 4 54. Wyres, K. *et al.* Genomic surveillance of antimicrobial resistant bacterial colonisation and  
5 infection in intensive care patients. *medRxiv* 2020.11.03.20224881 (2020)  
6 doi:10.1101/2020.11.03.20224881.
- 7 55. De Maio, N. *et al.* Comparison of long-read sequencing technologies in the hybrid  
8 assembly of complex bacterial genomes. *Microb. Genomics* **5**, (2019).
- 9 56. Decraene, V. *et al.* A Large, Refractory Nosocomial Outbreak of *Klebsiella pneumoniae*  
10 Carbapenemase-Producing *Escherichia coli* Demonstrates Carbapenemase Gene  
11 Outbreaks Involving Sink Sites Require Novel Approaches to Infection Control.  
12 *Antimicrob. Agents Chemother.* **62**, (2018).
- 13 57. Greig, D., Dallman, T. & Jenkins, C. Oxford Nanopore sequencing elucidates a novel  
14 stx2f carrying prophage in a Shiga toxin producing *Escherichia coli*(STEC) O63:H6  
15 associated with a case of haemolytic uremic syndrome (HUS). *Access Microbiol.* **1**, 782  
16 (2019).
- 17 58. Greig, D. R., Dallman, T. J., Hopkins, K. L. & Jenkins, C. MinION nanopore sequencing  
18 identifies the position and structure of bacterial antibiotic resistance determinants in a  
19 multidrug-resistant strain of enteroaggregative *Escherichia coli*. *Microb. Genomics* **4**,  
20 e000213 (2018).
- 21 59. Rang, F. J., Kloosterman, W. P. & de Ridder, J. From squiggle to basepair: computational  
22 approaches for improving nanopore sequencing read accuracy. *Genome Biol.* **19**, 90  
23 (2018).
- 24 60. Kato, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple  
25 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066  
26 (2002).

- 1 61. Treangen, T. J., Ondov, B. D., Koren, S. & Phillippy, A. M. The Harvest suite for rapid  
2 core-genome alignment and visualization of thousands of intraspecific microbial  
3 genomes. *Genome Biol.* **15**, 524 (2014).
- 4 62. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
5 large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
- 6 63. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant  
7 bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
- 8 64. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using  
9 MinHash. *Genome Biol.* **17**, 132 (2016).
- 10 65. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial  
11 genome assemblies from short and long sequencing reads. *PLOS Comput. Biol.* **13**,  
12 e1005595 (2017).
- 13 66. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome  
14 assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
- 15 67. Walker, B. J. *et al.* Pilon: An Integrated Tool for Comprehensive Microbial Variant  
16 Detection and Genome Assembly Improvement. *PLOS ONE* **9**, e112963 (2014).
- 17 68. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.*  
18 **5**, R12 (2004).
- 19 69. Köster, J. & Rahmann, S. Snakemake—a scalable bioinformatics workflow engine.  
20 *Bioinforma. Oxf. Engl.* **34**, 3600 (2018).
- 21 70. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
22 *ArXiv13033997 Q-Bio* (2013).
- 23 71. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat.*  
24 *Methods* **9**, 357–359 (2012).
- 25 72. Merkel, D. Docker: lightweight Linux containers for consistent development and  
26 deployment. *Linux J.* **2014**, 2:2 (2014).

- 1 73. Kurtzer, G. M., Sochat, V. & Bauer, M. W. Singularity: Scientific containers for mobility of
- 2 compute. *PLOS ONE* **12**, e0177459 (2017).