

1 **StrainFLAIR: Strain-level profiling of** 2 **metagenomic samples using variation** 3 **graphs**

4 **Kévin Da Silva^{1,2,*}, Nicolas Pons², Magali Berland², Florian Plaza Oñate²,**
5 **Mathieu Almeida², and Pierre Peterlongo¹**

6 ¹**Univ Rennes, Inria, CNRS, IRISA - UMR 6074, F-35000 Rennes, France**

7 ²**Université Paris-Saclay, INRAE, MGP, 78350 Jouy-en-Josas, France**

8 Corresponding author:

9 *Kévin Da Silva kevin.da-silva@inria.fr

10 Email address:

11 **ABSTRACT**

12 Current studies are shifting from the use of single linear references to representation of multiple genomes
13 organised in pangenome graphs or variation graphs. Meanwhile, in metagenomic samples, resolving
14 strain-level abundances is a major step in microbiome studies, as associations between strain variants
15 and phenotype are of great interest for diagnostic and therapeutic purposes.

16 We developed *StrainFLAIR* with the aim of showing the feasibility of using variation graphs for indexing
17 highly similar genomic sequences up to the strain level, and for characterizing a set of unknown sequenced
18 genomes by querying this graph.

19 On simulated data composed of mixtures of strains from the same bacterial species *Escherichia coli*,
20 results show that *StrainFLAIR* was able to distinguish and estimate the abundances of close strains, as
21 well as to highlight the presence of a new strain close to a referenced one and to estimate its abundance.
22 On a real dataset composed of a mix of several bacterial species and several strains for the same species,
23 results show that in a more complex configuration *StrainFLAIR* correctly estimates the abundance of
24 each strain. Hence, results demonstrated how graph representation of multiple close genomes can be
25 used as a reference to characterize a sample at the strain level.

26 **Availability:** <http://github.com/kevsilva/StrainFLAIR>

27 **INTRODUCTION**

28 The use of reference genomes has shaped the way genomics studies are currently conducted. Reference
29 genomes are particularly useful for reference guided genomic assembly, variant calling or mapping
30 sequencing reads. For the later, they provide a unique coordinate system to locate variants, allowing
31 to work on the same reference and easily share information. However, the usage of reference genomes
32 represented as flat sequences reaches some limits (Ballouz et al., 2019).

33 Close reference genomes or genomes of strains from the same species show a high sequence similarity.
34 Mapping sequencing reads on similar reference genomes results in mis-mapped reads or ambiguous
35 alignments generating noise in the downstream analysis, that has yet to be clarified (Na et al., 2016). This
36 has led recent methods to provide a representation of multiple genomes as genome graphs, also called
37 variation graphs, in which each path is a different known variation. Such graph representations are well
38 defined, and tools to build and manipulate graphs are under active development (Garrison et al., 2017;
39 Kim et al., 2019; Rakocevic et al., 2019; Li et al., 2020).

40 This graph structure provides obvious advantages such as the reduction of the data redundancy, while
41 highlighting variations (Garrison et al., 2018). However, it also introduces novel difficulties. Updating
42 a graph with novel sequences, adapting existing efficient algorithms for read mapping, and, mainly,
43 developing new ways to analyse sequence-to-graph mapping results for downstream analyses are among
44 those new challenges. The work presented here primarily focuses on this latest point and proposes to
45 show the feasibility of using a variation graph for identifying and estimating abundances, at the strain

46 level, from an unknown metagenomic read set.

47 In the context of metagenomics, representing genomes in graphs is of particular interest for indexing
48 microorganism genomes. Microorganisms are predominant in almost every ecosystems from ocean
49 water (Sunagawa et al., 2015) to human body (Clemente et al., 2012), and play major functioning roles
50 in them (New and Brito, 2020). While studies in microbial ecology are facing a bottleneck due to the
51 difficulty of isolating and cultivating most of those microbes in laboratory, preventing the analysis of
52 the complex structure and dynamics of the microbial communities (Stewart, 2012), high-throughput
53 sequencing in metagenomics offers the opportunity to study a whole ecosystem. In particular, shotgun
54 sequencing allows a resolution up to the species level (Jovel et al., 2016), and enable samples analysis
55 in terms of population stratification, microbial diversity or bio-markers identification (Quince et al., 2017).
56 Understanding of microbial communities structure and dynamics is usually revealed by resolving the
57 species present in samples and their relative abundances, which can then be associated with phenotypes,
58 notably in the field of human health (Ehrlich, 2011; Vieira-Silva et al., 2020; Solé et al., 2021). Now,
59 characterizing samples at the strain level has a growing interest, as it may highlight new associations with
60 phenotypes, and a better understanding of the functional impact of strains in host-microbe interactions
61 is crucial to new therapeutic strategies and personalized medicine. *Escherichia coli*, which has a highly
62 variable genome, is a well-known example since some strains are harmless commensals in the human
63 gut microbiota while others are harmful pathogens (Rasko et al., 2008; Loman et al., 2013). Current
64 approaches to handle multiple similar genomes as with strains use gene clustering and then select the
65 representative sequence of each cluster, getting rid of the redundancy but also the variations, yet crucial
66 to distinguish the strains of a species (Qin et al., 2010). Hence, indexation of a set of known strains is a
67 good framework for testing the ability of a variation graph to capture the diversity while offering a way to
68 correctly assign sequenced data to the strains they belong to.

69 In this work, we present `StrainFLAIR`, a novel method and its implementation that uses variation
70 graph representation of gene sequences for strain identification and quantification. We proposed novel
71 algorithmic and statistical solutions for managing ambiguous alignments and computing an adequate
72 abundance metric at the graph node level. Results have shown that we could correctly identify and quantify
73 strains present in a sample. Notably, we could also identify close strains not present in the reference.

74 `StrainFLAIR` is available at <http://github.com/kevsilva/StrainFLAIR>.

75 METHODS

76 We propose here a description of our tool `StrainFLAIR` (STRAIN-level proFiLing using vArLation
77 gRaph). This method exploits various state-of-the-art tools and proposes novel algorithmic solutions
78 for indexing bacterial genomes at the strain-level. It also permits to query metagenomes for assessing
79 and quantifying their content, in regards to the indexed genomes. An overview of the index and query
80 pipelines are presented on Fig. 1.

81 Rational for the choice of third-party tools and their detailed usages are given in Supplementary
82 Materials, Section S1.1.

83 Indexing strains

84 *Gene prediction*

85 As non-coding DNA represents 15% in average of bacterial genomes and is not well characterized in
86 terms of structure, `StrainFLAIR` focuses on protein-coding genes in order to characterize strains by
87 their gene content and nucleotidic variations of them. Moreover, non-coding DNA regions can be highly
88 variable (Thorpe et al., 2017) and taking into account complete genomes would then lead to highly
89 complex graphs, and combinatorial explosions when mapping reads. Additionally, complete genomes
90 are not always available. Focusing on the genes allows to use also drafts and metagenome-assembled
91 genomes or a pre-existing set of known genes (Qin et al., 2010; Li et al., 2014). Hence, `StrainFLAIR`
92 indexes genes instead of complete genomes in graphs.

93 Genes are predicted using `Prodigal`, a tool for prokaryotic protein-coding genes prediction (Hyatt
94 et al., 2010).

95 Knowing that some reads map at the junction between the gene and intergenic regions, by conserving
96 only gene sequences, mapping results are biased towards deletions and drastically lower the mapping
97 score. In order to alleviate this situation, we extend the predicted gene sequences at both ends. Hence,

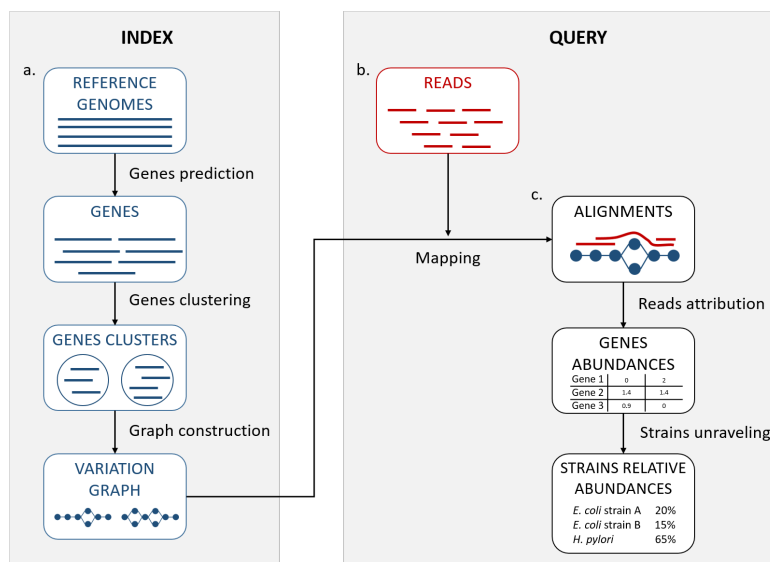


Figure 1. StrainFLAIR overview. a. Indexation. Input is a set of known reference genomes of various bacterial species and strains. StrainFLAIR uses a graph for indexing genes of those reference genomes. **b. Read mapping** on the previously mentioned graph. **c. Mapped reads analysis.** StrainFLAIR assigns and estimates species and strain abundances of a bacterial metagenomic sample represented as short reads.

98 StrainFLAIR conserves predicted genes plus their surrounding sequences. By default, and if the
 99 sequence is long enough, we conserve 75 bp on the left and on the right of each gene.

100 **Gene clustering**

101 Genes are clustered into gene families using CD-HIT (Li and Godzik, 2006). For the clustering step, the
 102 genes without extensions are used in order to strictly cluster according to the exact gene sequences and
 103 no parts of intergenic regions. CD-HIT-EST is used to realize the clustering with an identity threshold
 104 of 0.95 and a coverage of 0.90 on the shorter sequence. The local sequence identity is calculated as the
 105 number of identical bases in alignment divided by the length of the alignment. Sequences are assigned to
 106 the best fitting cluster verifying these requirements.

107 **Graph construction**

108 Each gene family is represented as a variation graph (Fig. 2). Variation graphs are bidirected DNA
 109 sequence graphs that represents multiple sequences, including their genetic variation. Each node of the
 110 graph contains sub-sequences of the input sequences, and successive nodes draw paths on the graph.
 111 Paths corresponding to reference sequences are specifically called “colored paths”. Each colored path
 112 corresponds to the original sequences of a gene in the cluster.

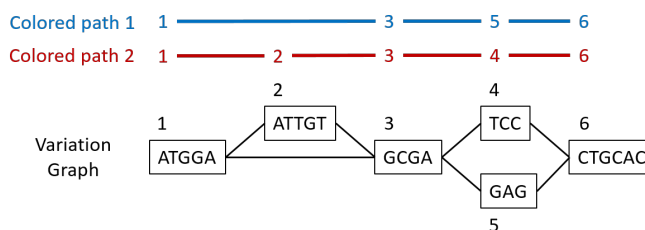


Figure 2. Illustration of a variation graph structure and colored paths. Each node of the graph contains a sub-sequence of the input sequences and is integer-indexed. A path corresponding to an input sequence is called a colored path, and is encoded by its succession of node ids, e.g. 1,3,5,6 for the colored path 1 in this example.

113 In the case of a cluster composed of only one sequence, `vg toolkit` (Garrison et al., 2017)
114 is used to convert the sequence into a flat graph. Alternatively, when a cluster is composed of two
115 sequences or more, `minimap2` (Li, 2018) is used to generate a multiple sequence alignment. Then
116 `seqwish` (Garrison, 2021) is used to convert this multiple sequence alignment into a variation graph.
117 All the so-computed graphs (one per input cluster) are then concatenated to produce a single variation
118 graph where each cluster of genes is a connected component.

119 The index is created once for a set of reference genomes. Afterward, any set of sequenced reads can
120 be profiled at the strain-level based on this index.

121 **Querying variation graphs**

122 ***Mapping reads***

123 For mapping reads on the previously described reference graph, we use the sequence-to-graph mapper `vg`
124 `mpmap` from `vg toolkit`. It produces a so-called “multipath alignments”. A multipath alignment is a
125 graph of partial alignments and can be seen as a sub-graph (a subset of edges and vertices) of the whole
126 variation graph (see Fig. 3 for an example). The mapping result describes, for each read, the nodes of the
127 variation graph traversed by the alignment and the potential mismatches or indels between the read and
128 the sequence of each traversed node.

129 ***Reads attribution***

130 When mapping a read on a graph with colored path, two key issues arise, as illustrated Fig. 3. As mapping
131 generates a sub-graph per mapped read, the most probable mapped path(s) has / have to be defined. In the
132 meanwhile, the most probable mapped path(s) corresponding to a colored path also have to be defined.
133 Hence we developed an algorithm to analyse and convert, when possible, a mapping result into one or
134 several continuous path(s) (successive nodes joined by only one edge) per mapped read. In addition we
135 propose an algorithm to attribute such path to most probable colored path(s).

136 ***Path attribution***

137 A breadth first search on the multipath alignment is proposed. It starts at each node of the alignment
138 with a user-defined threshold on the mapping score. A single path alignment with a mapping score
139 below this threshold is ignored, and the single path alignment with the best mapping score is retained.
140 Additionally, for each alignment, nodes are associated with a so-called “horizontal coverage” value. The
141 horizontal coverage of a node by a read corresponds to the proportion of bases of the node covered by the
142 read. Hence, a node has an horizontal coverage of 1 if all its nucleotides are covered by the read with or
143 without mismatches or indels.

144 Because of possible ties in mapping score, the search can result in multiple single path alignments, as
145 illustrated Fig. 3(A). This situation corresponds to a read which sequence is found in several different
146 genes or to a read mapping onto the similar region of different versions of a gene.

147 To take into account ambiguous mapping affectations, as shown below, the parsing of the mapping
148 output is decomposed into two steps. The first step processes the reads that mapped only a unique colored
149 path (called “unique mapped reads” here), corresponding to a single gene. The second step processes the
150 reads with multiple alignments (called “multiple mapped reads” here).

151 ***Colored path attribution***

152 Once a read is assigned to one or several path alignment(s), it still has to be attributed, if possible, to a
153 colored path. The following process attributes each mapped read to a colored path and various metrics for
154 downstream analyses are computed. In particular, an absolute abundance for each node of the variation
155 graph, called the “node abundance”, is computed, first focusing on unique mapped reads (first step). For a
156 given alignment, the successive nodes composing the path are compared to the existing colored paths of
157 the variation graph. If the alignment matches part of a colored path, the number of mapped reads on this
158 path is incremented by one (i.e. reads raw count). The node abundance for each node of the alignment is
159 incremented with its horizontal node coverage defined by this alignment. Alignments with no matching
160 colored paths are skipped.

161 Then, we focus on multiple mapped reads (second step), as illustrated Fig. 3(B). During this step, the
162 alignment matches multiple colored paths. Hence, the abundance is distributed to each matching colored
163 path relatively to the ratio between them. This ratio is determined from the reads raw count of each path
164 from the first step. For example, if 70 unique mapped reads were found for path1 and 30 for path2 during
165 the first step, a read matching ambiguously both path1 and path2 during the second step counts as 0.7 for

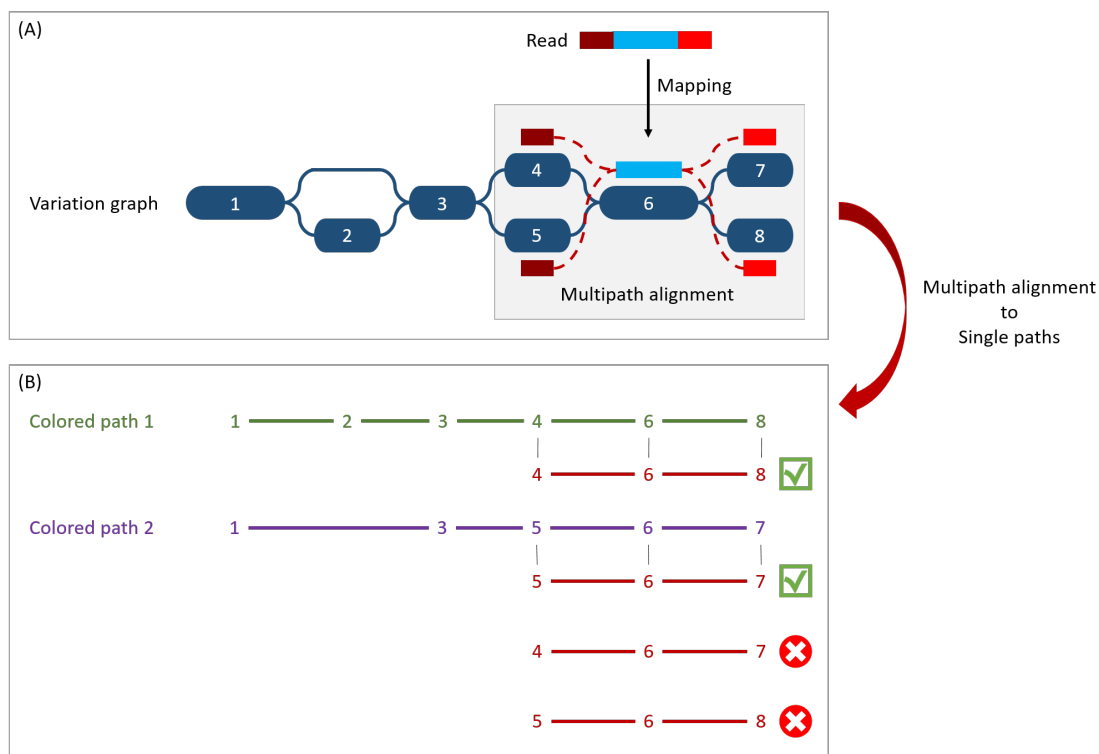


Figure 3. Illustration of the multipath alignment concept and the read attribution process. (A) **Path attribution.** The region of the read in blue aligns un-ambiguously to a node of the graph while the dark and light red parts can either align to the top or the bottom nodes of their respective mapping localization (due to mismatches that can align on both nodes for example), drawing an alignment as a sub-graph of the reference variation graph, and thus opening the possibility of four single path alignments. (B) **Colored path attribution.** First, from the multipath alignment (all four read sub-paths), the breadth search finds the possible corresponding single path alignments while respecting the mapping score threshold imposed by the user. Here, for the example, all four possible paths are considered valid. Second, each single path is compared to the colored paths from the reference variation graph. Two single path alignments matched the colored paths (4-6-8 and 5-6-7). As it mapped equally more than one colored path, this read falls in the multiple mapped reads case and is processed during the second step of the algorithm.

166 path1 and 0.3 for path2. This ratio is applied to increment both the raw count of reads and the coverage of
167 the nodes.

168 **Gene-level and strain-level abundances**

169 *StrainFLAIR* output is decomposed into an intermediate result describing the queried sample and
170 gene-level abundances, and the final result describing the strain-level abundances.

171 **Gene-level**

172 After parsing the mapping result, the first output provides information for each colored path, *i.e.*
173 each version of a gene. Thereby, this first result proposes gene-level information including abundances.
174 Exhaustive description of these intermediate results is provided in Section S1.2 in Supplementary Materials.
175 We describe here three major metrics outputted by *StrainFLAIR*:

176 **The mean abundance of the nodes composing the path.** Instead of solely counting reads, we make
177 full use of the graph structure and we propose abundances computation for each node as previously
178 explained, and as already done for haplotype resolution (Baaijens et al., 2019). Hence, for each colored
179 path, the gene abundance is estimated by the mean of the nodes abundance.

180 In order to not underestimate the abundance in case of a lack of sequencing depth (which could result
181 in certain nodes not to be traversed by sequencing reads), the **mean abundance without the nodes of**

182 **the path never covered by a read** is also outputted.

183 The mean abundance with and without these non-covered nodes are computed using unique mapped
184 reads only or all mapped reads.

185 The **ratio of covered nodes**, defined as the proportion of nodes from the path which abundance is
186 strictly greater than zero.

187 **Strain-level**

188 Strain-level abundances are then obtained by exploiting the specific genes of each reference genome
189 from these intermediate results. First, for each genome, the proportion of detected genes is computed,
190 as the proportion of specific genes on which at least one read maps. Then, the global abundance of the
191 genome is computed as the mean or median of all its specific gene abundances. However, if the proportion
192 of detected genes is less than a user-defined threshold, the genome is considered absent and hence its
193 abundance is set to zero.

194 *StrainFLAIR* final output is a table where each line corresponds to one of the reference genomes,
195 containing in columns the proportion of detected specific genes, and our proposed metrics to estimate their
196 abundances (using mean or median, with or without never covered nodes as described for the gene-level
197 result).

198 Results presented Section S1.3 in Supplementary Materials validate and motivate the proposed
199 abundance metric by comparing it to the expected abundances and other estimations using linear models.

200 **RESULTS**

201 We validated our method on both a simulated and a real dataset. All computations were performed using
202 *StrainFLAIR*, version 0.0.1, with default parameters. The relative abundances estimation was based
203 on the mean of the specific gene abundances, computed by taking into account all the nodes (including
204 non-covered nodes), and using a threshold on the proportion of detected specific genes of 50%.

205 Results were compared to *Kraken2* (Wood et al., 2019) considered as one of the state-of-the-art tool
206 dedicated to the characterization of read set content, and based on flat sequences as references. Read
207 counts given by *Kraken2* were normalized by the genome length and converted into relative abundances.

208 Computing setup and performances are indicated in Supplementary Materials, Section S1.4.

209 **Validation on a simulated dataset**

210 We first validated our method on simulated data, focusing on a single species with multiple strains. Our
211 aim was to validate the *StrainFLAIR* ability to identify and quantify strains given sequencing data
212 from a mixture of several strains of uneven abundances, and with one of them absent from the index.

213 **Reference variation graph**

214 We selected complete genomes of *Escherichia coli*, a predominant aerobic bacterium in the gut micro-
215 biota (Tenaillon et al., 2010), and a species known for its phenotypic diversity (pathogenicity, antibiotics
216 resistance) mostly resulting from its high genomic variability (Dobrindt, 2005).

217 Eight strains of *E. coli* were selected for this experiment from the NCBI¹. Seven were used to construct
218 a variation graph (*E. coli* IAI39, O104:H4 str. 2011C-3493, str. K-12 substr. MG1655, SE15, O157:H16
219 str. Santai, O157:H7 str. Sakai, O26 str. RM8426), and one was used as an unknown strain in a strains
220 mixture (*E. coli* BL21-DE3).

221 **Mixtures and sequencing simulations**

222 Our aim was to simulate the co-presence of several *E. coli* strains. Two simulations with sequencing
223 errors were conducted in order to highlight the detection and quantification of strains in a mixture. For
224 each one, we tested our approach with various read coverage, as described below.

225 We simulated the sequencing of three strains to mimic complex single species composition in
226 metagenomic samples. One of the strain was in equal abundance of one of the two others, potentially
227 making it more difficult to distinguish, or in lower abundance, potentially making it more difficult to
228 detect at all. The first simulation was a mixture composed of three strains contributing in the reference
229 graph: *E. coli* O104:H4 2011c-3493, IAI39, and K-12 MG1655. The second simulation was a mixture
230 composed of three strains: *E. coli* O104:H4 2011c-3493, IAI39, and BL21-DE3. The later being absent
231 from the reference variation graph thus simulating a new strain to be identified and quantified.

¹[https://www.ncbi.nlm.nih.gov/genome/?term=txid562\[orgn\]](https://www.ncbi.nlm.nih.gov/genome/?term=txid562[orgn])

232 For both simulations, short sequencing reads of 150 bp were simulated using `vg sim` from `vg`
 233 `toolkit` with a probability of errors set to 0.1% : 300,000 reads for *E. coli* O104:H4 2011c-3493
 234 (representing $\approx 8.5x$), 200,000 reads for *E. coli* IAI39 (representing $\approx 5.8x$). For both simulations, various
 235 quantities of reads were generated for K-12 MG1655 or BL21-DE3: 200,000, 100,000, 50,000, 25,000,
 236 10,000, 5,000 or 1,000 reads, representing approximately 6.5x, 3x, 1.6x, 0.8x, 0.3x, 0.2x, and 0.03x
 237 respectively for these two strains.

238 **Strain-level abundances**

239 As explained in Methods, we computed the strain-level abundances using the specific gene-level abundance
 240 table obtained by mapping the simulated reads onto the variation graph. We compared our results to the
 241 expected simulated relative abundances.

#reads K-12	Method	O104:H4	IAI39	K-12	Sakai	SE15	Santai	RM8426
1,000	Expected	59.88	39.92	0.2	0	0	0	0
	StrainFLAIR	56.45	43.55	0	0	0	0	0
	Kraken2	38.91	60.72	0.22	0.04	0.07	0.03	0.02
25,000	Expected	57.14	38.1	4.76	0	0	0	0
	StrainFLAIR	52.1	40.58	7.32	0	0	0	0
	Kraken2	37.23	58.1	4.51	0.04	0.07	0.03	0.02
200,000	Expected	42.86	28.57	28.57	0	0	0	0
	StrainFLAIR	38.12	29.83	32.05	0	0	0	0
	Kraken2	28.31	44.18	27.35	0.04	0.08	0.03	0.02

Table 1. Reference strains relative abundances expected and computed by StrainFLAIR or Kraken2 for each simulated experiment with variable coverage of the K-12 MG1655 strain. Best results are shown in bold. Complete results are presented Section S1.6 in Supplementary Materials.

242 **Simulation 1: mixtures with K-12 MG1655, present in the reference graph**

243 StrainFLAIR successfully estimated the relative abundances of the three strains present in the
 244 mixture (Table 1), the sum of squared errors between the estimation given by our tool and the expected
 245 relative abundance was between 25 and 45 for all the experiments. However, it did not detect the very
 246 low abundant strain in the case of the mixture with 1,000 simulated reads for K-12 MG1655 (coverage of
 247 $\approx 0.03x$). With our methodology, the threshold on the proportion of detected genes (see Methods) lead
 248 to set relative abundance to zero of likely absent strains. This reduces both the underestimation of the
 249 relative abundances of the present strains and the overestimation of the absent strains.

250 In comparison, Kraken2 did not provide this resolution. Applied to our simulated mixtures, while
 251 Kraken2 was slightly better for K-12 MG1655 abundance estimation, it overestimated IAI39 relative
 252 abundance and underestimated O104's one, leading to an overall higher sum of squared errors (between
 253 456 and 872) compared to the expected abundances. Moreover, it set relative abundances to all the seven
 254 reference strains whereas four of them were absent from the mixture. This was expected as some reads
 255 (from intergenic regions for example) can randomly be similar to regions of genes from absent strains.

256 **Simulation 2: mixtures with BL21-DE3, absent from the reference graph**

257 Here, BL21-DE3 was considered an unknown strain, not contributing to the variation graph. The closest
 258 strain of BL21-DE3 in the graph, according to fastANI (Jain et al., 2018), was K-12 MG1655 (98.9%
 259 of identity, see Supplementary Materials, Section S1.5). Thus we expected to find signal of BL21-DE3
 260 through the results on K-12 MG1655.

261 As with the K-12 MG1655 mixtures, StrainFLAIR successfully estimated the relative abundances
 262 of the two known strains present in the mixture (Table 2), the sum of squared errors between the estimation
 263 given by our tool and the expected relative abundance was between 22 and 180 for all the experiments.
 264 Labelled as K-12, it also gave close estimations for BL21-DE3. Again, it did not detect the very low
 265 abundant strain in the case of the mixture with 1,000, 5,000, and 10,000 simulated reads for BL21-DE3.
 266 Also similarly to the K-12 MG1655 mixtures experiments, Kraken2 overestimated IAI39 relative
 267 abundance and underestimated O104's one (sum of squared errors between 751 and 873), even less

#reads BL21-DE3	Method	O104:H4	IAI39	K-12	Sakai	SE15	Santai	RM8426
1,000	Expected	59.88	39.92	(0.2)	0	0	0	0
	StrainFLAIR	56.47	43.53	0	0	0	0	0
	Kraken2	38.93	60.76	0.11	0.05	0.08	0.04	0.03
25,000	Expected	57.14	38.1	(4.76)	0	0	0	0
	StrainFLAIR	54.09	41.71	4.2	0	0	0	0
	Kraken2	37.75	58.93	2.16	0.28	0.34	0.25	0.29
200,000	Expected	42.86	28.57	(28.57)	0	0	0	0
	StrainFLAIR	46.95	35.34	17.72	0	0	0	0
	Kraken2	31.14	48.83	13.53	1.57	1.67	1.58	1.68

Table 2. Reference strain relative abundances expected and computed by StrainFLAIR or Kraken2 for each simulated experiment with variable coverage of the BL21-DE3 strain, absent from the reference variation graph. BL21-DE3 strain expected abundances are given in parentheses in the K-12 column. Best results are shown in bold. Complete results are presented Section S1.6 in Supplementary Materials.

268 precisely than in the previous experiment. With sufficient coverage (here from the 0.8x for BL21-DE3),
269 StrainFLAIR was closer to the expected values for all the reference strains than Kraken2.

270 Interestingly, the proportion of detected specific genes for each strain (Fig. 4) seems to highlight a
271 pattern allowing to distinguish present strains, absent strains and likely new strains close to the reference
272 in the graph. According to the experiments with enough coverage (from 25,000 simulated reads for
273 BL21-DE3), three groups of proportions could be observed: proportion of almost 100% (O104:H4 and
274 IAI39 : strains present in the mixtures and in the reference graph), proportion under 30-35% (Sakai, SE15,
275 Santai, and RM8426 : strains absent from the mixtures), and an in-between proportion around 60-70% for
276 K-12 MG1655 (closest strain to BL21-DE3).

277 It was expected that an absent strain would have specific genes detected as StrainFLAIR detects a
278 gene once only one read mapped on it. However, all absent strains had a proportion at around 30% except
279 K-12 MG1655 which proportion was twice higher. Conjointly with the non-null abundance estimated for
280 the reference K-12 MG1655, this suggests the presence of a new strain whose genome is highly similar to
281 K-12 MG1655.

282 Validation on a real dataset

283 We used a mock dataset available on EBI-ENA repository under accession number PRJEB42498, in order
284 to validate our method on real sequencing data from samples composed of various species and strains.
285 The mock dataset is composed of 91 strains of bacterial species for which complete genomes or sets of
286 contigs are available, including plasmids. Among the species, two of them contained each two different
287 strains. Three mixes had been generated from the mock, and we used the “Mix1A” in the following
288 results.

289 Even though 20 out of 91 strains were absents in this mix, we indexed the full set of 91 genomes.
290 This was done in order to mimic a classical StrainFLAIR use case where the queried data is mainly
291 unknown, and the reference graph contains species or strains not existing in these queried data. The
292 metagenomic sample was sequenced using Illumina HiSeq 3000 technology and resulted in 21,389,196
293 short paired-end reads.

294 We compared our results to the expected abundances of each strain in the sample defined as the
295 theoretical experimental DNA concentration proportion. As such, it has to be noted that potential
296 contamination and/or experimental bias could have occurred and affected the expected abundances.

297 Strain detection

298 Among the 91 strains used in the reference variation graph, StrainFLAIR detected 65 strains. All of
299 these 65 strains were indeed sequenced in Mix1A. Hence, StrainFLAIR produced no false positive.
300 From the 26 strains considered absent by StrainFLAIR, 20 were not present in the sample (true
301 negatives) and 6 should have been detected (false negatives). However, the term false negative has to be
302 soften as the ground truth remains uncertain. Among those 6 undetected strains, all of them had theoretical
303 abundance below 0.1%.

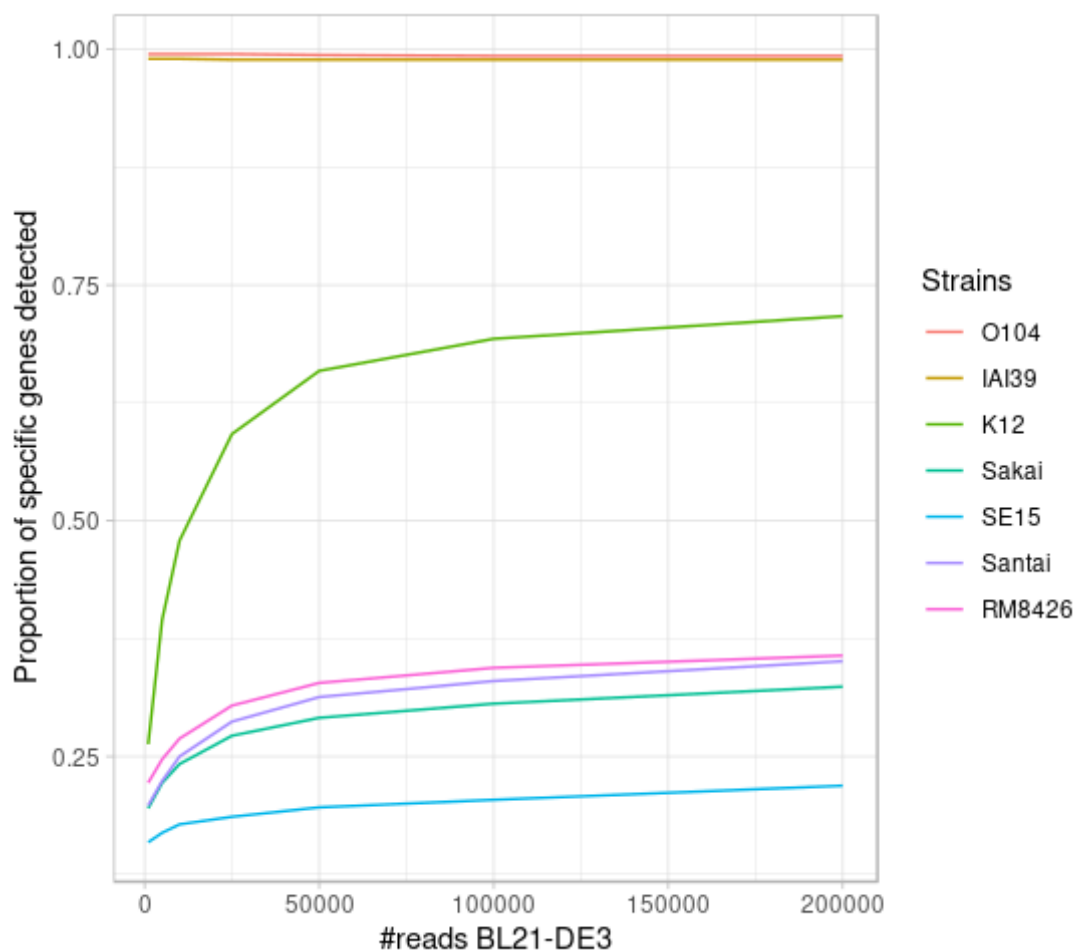


Figure 4. Proportion of detected specific genes for each simulated experiment with variable coverage of the BL21-DE3 strain, absent from the reference graph.

304 More precisely, among the 6 strains undetected by *StrainFLAIR*, 5 had some detected genes,
305 but below the 50% threshold. In this case, by default, *StrainFLAIR* discards these strains. Finally,
306 only one of the undetected strains (*Desulfovibrio desulfuricans* ND 132) should have been theoretically
307 detected (even if its expected coverage was below 0.1%), but no specific gene was identified. Considering
308 that *StrainFLAIR* uses a permissive definition of detected gene (at least one read maps on the gene),
309 having strictly no specific genes detected for *Desulfovibrio desulfuricans* ND 132 suggests that this strain
310 might in fact be absent from Mix1A. This is also supported by the result from *Kraken2* which estimated
311 a relative abundance of $\approx 9e-5$, almost 500 times lower than the theoretical result.

312 As in the simulated dataset validation, *Kraken2* affected non-null abundances to all the references
313 and thus could not be used to definitely conclude on presence/absence of strains in the sample.

314 **Strain relative abundances**

315 For the estimated relative abundances, *StrainFLAIR* gave more similar results compared to the
316 state-of-the-art tool *Kraken2* than the experimental values (Fig. 5). The sum of squared error between
317 *StrainFLAIR* and *Kraken2* was around 11. *StrainFLAIR* and *Kraken2* gave similar results
318 compared to the experimental values, with sum of squared errors of around 209 and 211 respectively.

319 Interestingly, *Thermotoga petrophila* RKU-1 is the only case where results from *StrainFLAIR*
320 and *Kraken2* differs greatly, with, in addition, the theoretical abundance being in-between. Moreover,
321 *Thermotoga* sp. RQ2 is the strain expected to be absent that *Kraken2* estimates with the highest relative
322 abundance among the other expected absent strains, and the only one exceeding the relative abundances
323 of two present strains. Considering the previous results on the simulated mixtures and that *Thermotoga*

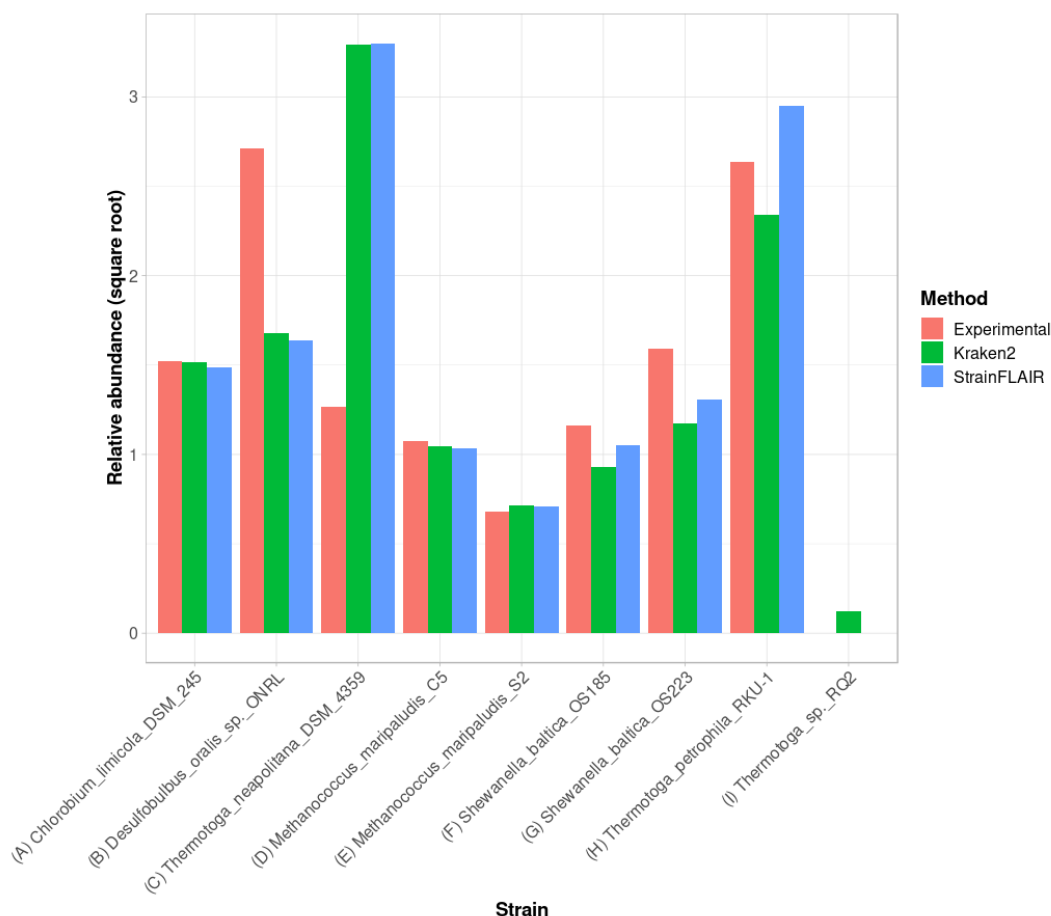


Figure 5. Experimental relative abundance compared to relative abundance as computed by StrainFLAIR and Kraken2. A selection of relevant results is shown here, see Supplementary Materials (Section S1.7) for the complete results. (A) Represents a case where StrainFLAIR and Kraken2 give similar results to the experimental value (18 cases over 91). (B) Represents a case where StrainFLAIR and Kraken2 give similar results, but lower than the experimental value (26 cases over 91). (C) Represents a case where StrainFLAIR and Kraken2 give similar results, but greater than the experimental value (16 cases over 91). (D, E, F, G) Represent the two species represented by two strains each. (H, I) Represent two atypical cases.

324 *petrophila* RKU-1 and *Thermotoga* sp. RQ2 are close species (fastANI around 96.6%) it could be an
 325 additional indicator of how tools like Kraken2 can be misled by too close species or strains.

326 In the sample, the species *Methanococcus maripaludis* was represented by two strains (S2 and C5) and
 327 the species *Shewanella baltica* likewise (OS223 and OS185). StrainFLAIR successfully distinguished
 328 and estimated the relative abundances of each strain of these two genomes. In this very situation and
 329 contrary to results on *E. coli* strains, Kraken2 was also able to correctly estimate the abundances.

330 DISCUSSION

331 Recent advances in sequencing technologies have provided large reference genome resources. Representa-
 332 tion and integration of those multiple genomes, often highly similar, are under active development and
 333 led to genome graphs based tools. Integrating multiple genomes from the same species is particularly
 334 interesting as it provides new opportunities to characterize strains, a key resolution, for instance opening
 335 the field of precision medicine (Albanese and Donati, 2017; Marchesi et al., 2016).

336 In this context, we developed StrainFLAIR, a new computational approach for strain level profiling
 337 of metagenomic samples, using variation graphs for representing all reference genomes. Our intention was

338 in the one hand to test whether or not indexing highly similar genomes in a graph enables to characterize
339 queried samples at the strain level, and, in the other hand, to provide a end-user tool able to perform the
340 indexation of genomes and the query of reads including the analyses of mapping results.

341 The method exploits state-of-the art-tools additionally to novel algorithmic and statistical solutions.
342 By indexing microbial species and/or strains in a graph, it enables the identification and quantification of
343 strains from a sequenced sample, mapped onto this graph.

344 We have demonstrated on simulated and on real datasets the ability of our method to identify and cor-
345 rectly estimate the abundance of microbial strains in metagenomic samples. In addition, *StrainFLAIR*
346 was able to highlight the presence and also to estimate a relative abundance for a strain similar to existing
347 references, but absent from these references.

348 We also showed that *StrainFLAIR* tended to set to zero the predicted abundance of low abundant
349 strains, while a tool like *Kraken2* was able detect them. As a result, it seemed that *StrainFLAIR*
350 loses the ability to detect very low abundant strains. However, in our simulations, this situation
351 corresponded to coverages of 0.03x or less, hence simulating a strain for which not all genomic content
352 was present. Eventually, it might be more relevant to define this strain as absent. Overall, there is a need to
353 distinguish between low abundant strains, insufficient sequencing depth, and reads from intergenic regions
354 or other genes randomly matching genes. In this regard, *StrainFLAIR* integrated a threshold on the
355 proportion of specific genes detected that can be further explored to refine which strain abundances are set
356 to zero. Importantly, results also showed that our graph-based tool had no false positive call, contrary to
357 general purpose tool *Kraken2* that detected 100% of strains that were indexed but absent from queried
358 reads.

359 From the validation on real datasets, we showed that *StrainFLAIR* was still able to correctly
360 estimate the relative abundances in a more complex context mixing both different species and different
361 strains, without being biased by references absent in the sample.

362 Our methodology taking into account all mapped reads and imposing a threshold that sets some strains
363 abundances to zero seems more adequate and closer to what is expected in reality. Moreover, being able
364 to detect some queried strains as absent is particularly interesting in the metagenomics context. Unlike
365 mock datasets that are of controlled and known compositions, no prior knowledge is available for real
366 metagenomic samples. They require the most exhaustive references - including unnecessary genomes -
367 hence strains absent from the sample. *StrainFLAIR* is a new step towards the objective to take into
368 account those unnecessary genomes without biasing the downstream analysis.

369 Measured computation time performances show that *StrainFLAIR* enables to analyse million reads
370 in a few hours. Even if this opens the doors to routine analyses of small read sets, new development
371 efforts will have to be made for reducing computation time in order to scale-up to very large datasets.

372 While *StrainFLAIR* focuses on profiling metagenomic samples at the strain level based on genes, it
373 opens the way to pangenomic studies. Genome graphs are used to capture all the information on variation
374 or similarity of sequences, which is particularly adapted to represent the gene repertoire diversity and the
375 set of nucleotidic variations found between the different genomes of a species. This work highlights the
376 importance to keep up working on pangenome graph representation.

377 The presence of queried unknown strain(s) is revealed both by reads mapping non-colored paths and
378 by the amount of nucleotidic variations (indels and substitutions). The natural continuation will be related
379 to the dynamical update of the graph when novel strains are detected in this way. This dynamicity will also
380 be particularly useful considering the future flow of new sequenced metagenomes and the development of
381 clinical metagenomics that will help to quickly and efficiently characterize in silico emerging strains of
382 human health interest.

383 **ACKNOWLEDGMENTS**

384 This work used the GenOuest bioinformatics core facility (<https://www.genouest.org>).

385 We acknowledge Mircea Podar for the providing of the mock dataset in premium access. Finally, we
386 thank Mahendra Mariadassou, Rayan Chikhi, Olivier Jaillon and David Vallenet for all their advice along
387 this work.

388 REFERENCES

- 389 Albanese, D. and Donati, C. (2017). Strain profiling and epidemiology of bacterial species from metage-
390 nomic sequencing. *Nature Communications*, 8(1):1–14.
- 391 Baaijens, J. A., der Roest, B. V., Köster, J., Stougie, L., and Schönhuth, A. (2019). Full-length de novo
392 viral quasispecies assembly through variation graph construction. *bioRxiv*, page 287177.
- 393 Ballouz, S., Dobin, A., and Gillis, J. (2019). Is it time to change the reference genome? *bioRxiv*, page
394 533166.
- 395 Clemente, J. C., Ursell, L. K., Parfrey, L. W., and Knight, R. (2012). The impact of the gut microbiota on
396 human health: An integrative view.
- 397 Dobrindt, U. (2005). (Patho-)Genomics of *Escherichia coli*.
- 398 Ehrlich, S. D. (2011). MetaHIT: The European Union project on metagenomics of the human intestinal
399 tract. In *Metagenomics of the Human Body*, pages 307–316. Springer New York.
- 400 Garrison, E. (2021). ekg/seqwish: alignment to variation graph inducer. [https://github.com/
401 ekg/seqwish](https://github.com/ekg/seqwish).
- 402 Garrison, E., Novak, A., Hickey, G., Eizenga, J., Dawson, E., Jones, W., Buske, O., and Lin, M. (2017).
403 Sequence variation aware references and read mapping with *vg*: the variation graph toolkit. *bioRxiv*.
- 404 Garrison, E., Sirén, J., Novak, A. M., Hickey, G., Eizenga, J. M., Dawson, E. T., Jones, W., Garg, S.,
405 Markello, C., Lin, M. F., Paten, B., and Durbin, R. (2018). Variation graph toolkit improves read
406 mapping by representing genetic variation in the reference.
- 407 Hyatt, D., Chen, G. L., LoCascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J. (2010). Prodigal:
408 Prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*, 11:119.
- 409 Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T., and Aluru, S. (2018). High throughput
410 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nature Communications*,
411 9(1):1–8.
- 412 Jovel, J., Patterson, J., Wang, W., Hotte, N., O’Keefe, S., Mitchel, T., Perry, T., Kao, D., Mason, A. L.,
413 Madsen, K. L., and Wong, G. K. (2016). Characterization of the gut microbiome using 16S or shotgun
414 metagenomics. *Frontiers in Microbiology*, 7(APR):459.
- 415 Kim, D., Paggi, J. M., Park, C., Bennett, C., and Salzberg, S. L. (2019). Graph-based genome alignment
416 and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology*, 37(8):907–915.
- 417 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–
418 3100.
- 419 Li, H., Feng, X., and Chu, C. (2020). The design and construction of reference pangenome graphs with
420 minigraph. *Genome Biology*, 21(1):265.
- 421 Li, J., Wang, J., Jia, H., Cai, X., Zhong, H., Feng, Q., Sunagawa, S., Arumugam, M., Kultima, J. R.,
422 Prifti, E., Nielsen, T., Juncker, A. S., Manichanh, C., Chen, B., Zhang, W., Levenez, F., Wang, J., Xu,
423 X., Xiao, L., Liang, S., Zhang, D., Zhang, Z., Chen, W., Zhao, H., Al-Aama, J. Y., Edris, S., Yang,
424 H., Wang, J., Hansen, T., Nielsen, H. B., Brunak, S., Kristiansen, K., Guarner, F., Pedersen, O., Doré,
425 J., Ehrlich, S. D., and Bork, P. (2014). An integrated catalog of reference genes in the human gut
426 microbiome. *Nature Biotechnology*, 32(8):834–841.
- 427 Li, W. and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or
428 nucleotide sequences. *Bioinformatics*, 22(13):1658–1659.
- 429 Loman, N. J., Constantinidou, C., Christner, M., Rohde, H., Chan, J. Z.-M., Quick, J., Weir, J. C., Quince,
430 C., Smith, G. P., Betley, J. R., Aepfelbacher, M., and Pallen, M. J. (2013). A Culture-Independent
431 Sequence-Based Metagenomics Approach to the Investigation of an Outbreak of Shiga-Toxigenic
432 *Escherichia coli* O104:H4. *JAMA*, 309(14):1502.
- 433 Marchesi, J. R., Adams, D. H., Fava, F., Hermes, G. D., Hirschfield, G. M., Hold, G., Quraishi, M. N.,
434 Kinross, J., Smidt, H., Tuohy, K. M., Thomas, L. V., Zoetendal, E. G., and Hart, A. (2016). The gut
435 microbiota and host health: A new clinical frontier. *Gut*, 65(2):330–339.
- 436 Na, J. C., Kim, H., Park, H., Lecroq, T., Léonard, M., Mouchard, L., and Park, K. (2016). FM-index of
437 alignment: A compressed index for similar strings. *Theoretical Computer Science*, 638:159–170.
- 438 New, F. N. and Brito, I. L. (2020). What Is Metagenomics Teaching Us, and What Is Missed?
- 439 Paten, B., Eizenga, J. M., Rosen, Y. M., Novak, A. M., Garrison, E., and Hickey, G. (2018). Superbubbles,
440 Ultrabubbles, and Cacti. In *Journal of Computational Biology*, volume 25, pages 649–663. Mary Ann
441 Liebert Inc.
- 442 Qin, J., Li, R., Raes, J., Arumugam, M., Burgdorf, K. S., Manichanh, C., Nielsen, T., Pons, N., Levenez,

- 443 F., Yamada, T., Mende, D. R., Li, J., Xu, J., Li, S., Li, D., Cao, J., Wang, B., Liang, H., Zheng, H., Xie,
444 Y., Tap, J., Lepage, P., Bertalan, M., Batto, J.-M., Hansen, T., Le Paslier, D., Linneberg, A., Nielsen,
445 H. B., Pelletier, E., Renault, P., Sicheritz-Ponten, T., Turner, K., Zhu, H., Yu, C., Li, S., Jian, M., Zhou,
446 Y., Li, Y., Zhang, X., Li, S., Qin, N., Yang, H., Wang, J., Brunak, S., Doré, J., Guarner, F., Kristiansen,
447 K., Pedersen, O., Parkhill, J., Weissenbach, J., MetaHIT Consortium, M., Bork, P., Ehrlich, S. D.,
448 and Wang, J. (2010). A human gut microbial gene catalogue established by metagenomic sequencing.
449 *Nature*, 464(7285):59–65.
- 450 Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J., and Segata, N. (2017). Shotgun metagenomics,
451 from sampling to analysis.
- 452 Rakocevic, G., Semenyuk, V., Lee, W. P., Spencer, J., Browning, J., Johnson, I. J., Arsenijevic, V., Nadj, J.,
453 Ghose, K., Suci, M. C., Ji, S. G., Demir, G., Li, L., Toptaş, B., Dolgoborodov, A., Pollex, B., Spulber,
454 I., Glotova, I., Kómar, P., Stachyra, A. L., Li, Y., Popovic, M., Källberg, M., Jain, A., and Kural, D.
455 (2019). Fast and accurate genomic analyses using genome graphs. *Nature Genetics*, 51(2):354–362.
- 456 Rasko, D. A., Rosovitz, M. J., Myers, G. S., Mongodin, E. F., Fricke, W. F., Gajer, P., Crabtree, J.,
457 Sebaihia, M., Thomson, N. R., Chaudhuri, R., Henderson, I. R., Sperandio, V., and Ravel, J. (2008).
458 The pangenome structure of *Escherichia coli*: Comparative genomic analysis of *E. coli* commensal and
459 pathogenic isolates. *Journal of Bacteriology*, 190(20):6881–6893.
- 460 Solé, C., Guilly, S., Da Silva, K., Llopis, M., Le-Chatelier, E., Huelin, P., Carol, M., Moreira, R.,
461 Fabrellas, N., De Prada, G., Napoleone, L., Graupera, I., Pose, E., Juanola, A., Borruel, N., Berland,
462 M., Toapanta, D., Casellas, F., Guarner, F., Doré, J., Solà, E., Ehrlich, S. D., and Ginès, P. (2021).
463 Alterations in Gut Microbiome in Cirrhosis as Assessed by Quantitative Metagenomics: Relationship
464 With Acute-on-Chronic Liver Failure and Prognosis. *Gastroenterology*, 160(1):206–218.e13.
- 465 Stewart, E. J. (2012). Growing unculturable bacteria.
- 466 Sunagawa, S., Coelho, L. P., Chaffron, S., Kultima, J. R., Labadie, K., Salazar, G., Djahanschiri, B., Zeller,
467 G., Mende, D. R., Alberti, A., Cornejo-Castillo, F. M., Costea, P. I., Cruaud, C., D’Ovidio, F., Engelen,
468 S., Ferrera, I., Gasol, J. M., Guidi, L., Hildebrand, F., Kokoszka, F., Lepoivre, C., Lima-Mendez, G.,
469 Poulain, J., Poulos, B. T., Royo-Llonch, M., Sarmiento, H., Vieira-Silva, S., Dimier, C., Picheral, M.,
470 Searson, S., Kandels-Lewis, S., Boss, E., Follows, M., Karp-Boss, L., Krzic, U., Reynaud, E. G., Sardet,
471 C., Sieracki, M., Velayoudon, D., Bowler, C., De Vargas, C., Gorsky, G., Grimsley, N., Hingamp, P.,
472 Iudicone, D., Jaillon, O., Not, F., Ogata, H., Pesant, S., Speich, S., Stemann, L., Sullivan, M. B.,
473 Weissenbach, J., Wincker, P., Karsenti, E., Raes, J., Acinas, S. G., and Bork, P. (2015). Structure and
474 function of the global ocean microbiome. *Science*, 348(6237).
- 475 Tenailon, O., Skurnik, D., Picard, B., and Denamur, E. (2010). The population genetics of commensal
476 *Escherichia coli*.
- 477 Thorpe, H. A., Bayliss, S. C., Hurst, L. D., and Feil, E. J. (2017). Comparative analyses of selection
478 operating on nontranslated intergenic regions of diverse bacterial species. *Genetics*, 206(1):363–376.
- 479 Vieira-Silva, S., Falony, G., Belda, E., Nielsen, T., Aron-Wisniewsky, J., Chakaroun, R., Forslund, S. K.,
480 Assmann, K., Valles-Colomer, M., Nguyen, T. T. D., Proost, S., Prifti, E., Tremaroli, V., Pons, N.,
481 Le Chatelier, E., Andreelli, F., Bastard, J. P., Coelho, L. P., Galleron, N., Hansen, T. H., Hulot, J. S.,
482 Lewinter, C., Pedersen, H. K., Quinquis, B., Rouault, C., Roume, H., Salem, J. E., Søndertoft, N. B.,
483 Touch, S., Alves, R., Amouyal, C., Galijatovic, E. A. A., Barthelemy, O., Batisse, J. P., Berland, M.,
484 Bittar, R., Blottière, H., Bosquet, F., Boubrit, R., Bourron, O., Camus, M., Cassuto, D., Ciangura,
485 C., Collet, J. P., Dao, M. C., Debedat, J., Djebbar, M., Doré, A., Engelbrechtsen, L., Fellahi, S.,
486 Fromentin, S., Giral, P., Graine, M., Hartemann, A., Hartmann, B., Helft, G., Herberg, S., Hornbak,
487 M., Isnard, R., Jaqueminet, S., Jørgensen, N. R., Julienne, H., Justesen, J., Kammer, J., Kerneis, M.,
488 Khemis, J., Krarup, N., Kuhn, M., Lampuré, A., Lejard, V., Levenez, F., Lucas-Martini, L., Massey,
489 R., Maziers, N., Medina-Stamminger, J., Moitinho-Silva, L., Montalescot, G., Moutel, S., Le Pavin,
490 L. P., Poitou-Bernert, C., Pousset, F., Pouzoulet, L., Schmidt, S., Silvain, J., Svendstrup, M., Swartz, T.,
491 Vanduyvenboden, T., Vatier, C., Verger, E., Walther, S., Dumas, M. E., Ehrlich, S. D., Galan, P., Götze,
492 J. P., Hansen, T., Holst, J. J., Køber, L., Letunic, I., Nielsen, J., Oppert, J. M., Stumvoll, M., Vestergaard,
493 H., Zucker, J. D., Bork, P., Pedersen, O., Bäckhed, F., Clément, K., and Raes, J. (2020). Statin therapy
494 is associated with lower prevalence of gut microbiota dysbiosis. *Nature*, 581(7808):310–315.
- 495 Wood, D. E., Lu, J., and Langmead, B. (2019). Improved metagenomic analysis with Kraken 2. *Genome
496 Biology*, 20(1):257.

497 **S1 SUPPLEMENTARY MATERIALS**

498 **S1.1 Third-party tools usage and rationale**

499 We propose here the motivations and precise usage of the third-party tools that are employed in
500 *StrainFLAIR*.

501 **S1.1.1 Graph construction**

502 *vg toolkit* allows to modify the graph including a normalization step. Normalization consists in
503 deleting redundant nodes (nodes containing the same sub-sequence and having the same parent and child
504 nodes), removing edges that do not introduce new paths, and merging nodes separated by only one edge.

505 For each cluster, if the colored paths of the corresponding graph still describe their respective input
506 sequences, the graph is normalized.

507 After the concatenation of all computed graphs (one for each cluster), the final single variation graph
508 is indexed using *vg toolkit*. Indexing a graph allows a fast querying of the graph when mapping
509 reads. Indexation uses two file formats: *XG*, which is a succinct graph index which presents a static
510 index of nodes, edges and paths of a variation graph, and *GCSA*, a generalized FM-index to directed
511 acyclic graphs. A *SNARLS* file is also generated, describing snarls (a generalization of the superbubble
512 concept (Paten et al., 2018)) in the variation graph and similarly allowing faster querying.

513 **S1.1.2 Mapping reads**

514 *vg toolkit* offers two sequence-to-graph mappers. The first one, *vg map*, outputs one or several
515 final paths for each alignment. However, in case of several alignments with equal mapping scores, only
516 one is randomly chosen. In order to get more exhaustive and accurate results, *StrainFLAIR* uses *vg*
517 *mpmap* to map reads on the variation graph.

518 The mapping results are given in *GAMP* format, then converted into *JSON* format with *vg toolkit*,
519 describing, for each read, the nodes of the graph traversed by the alignment.

520 **S1.2 Gene-level output by *StrainFLAIR***

521 Here we present the exhaustive description of information provided by *StrainFLAIR* at the gene level
522 (before strain-level computations). For each colored path *StrainFLAIR* provides the following items:

- 523 • The corresponding gene identifier.
- 524 • For each reference genome, the number of copies of the gene. Since each unique version of a gene
525 is represented once in the graph, whereas it can exist in several copies in the genome (duplicate
526 genes), the counts and abundances computed correspond to the sum of those copies. Keeping track
527 of the number of copies is important to normalize the counts.
- 528 • The cluster identifier to which the colored path belongs.
- 529 • For unique mapped reads: their raw number and their number normalized by the sequence length
530 (see Section Querying variation graphs in Methods).
- 531 • For unique plus multiple mapped reads: their raw number and their number normalized by the
532 sequence length (see Section Querying variation graphs in Methods).
- 533 • The mean abundance of the nodes composing the path, as defined in the manuscript.
- 534 • The mean abundance without the nodes of the path never covered by a read, as defined in the
535 manuscript.
- 536 • The ratio of covered nodes, as defined in the manuscript.

537 **S1.3 Abundance metrics validation**

538 The output of *StrainFLAIR* provides several metrics to estimate the abundance of the genes detected
539 in the sample.

540 For validation, we used a combination of LASSO (least absolute shrinkage and selection operator)
541 model and linear model on the simulated dataset to estimate the abundances at the strain-level, as the
542 abundance of a gene is a linear combination of the abundances of the strains it belongs to. As such,

543 we expect no intercept value for those models and have forced the intercept at zero for the following
544 modeling.

545 First, a LASSO model was used to perform strain selection. The response variable of the model was
546 the presence or absence of the genes according to the selected metric while the strains, described as their
547 genes content (number of copies), were the predictors. Then, a linear model was constructed with the
548 raw selected metric as the response variable, and only the strains selected by the LASSO model as the
549 predictors. The estimate of the strains relative abundance was thus the coefficients of the linear model
550 associated to the strains and transformed into relative values. For each metric, the sum of squared errors
551 between the real relative abundances and the estimated relative abundances from the linear model was
552 computed. The best metric was then defined as the one minimizing this sum of squared errors.

553 For the mixtures containing *E. coli* K-12 MG1655, the three expected strains were selected and thus
554 detected using LASSO, except for the mixture containing only 1,000 reads of K-12 MG1655 (representing
555 0.002% of the mixture, hence very negligible). For all the mixtures, the best metric was the mean
556 abundance computed from the node abundances and by taking into account the multiple mapped reads.

557 For the mixtures containing *E. coli* BL21-DE3, BL21-DE3 being absent from the reference but very
558 close to K-12 MG1655, we expected to get some detection of K-12 in the results. The three expected
559 strains were selected and thus detected using LASSO, except for the mixture containing only 1,000 reads
560 of BL21-DE3 (representing 0.002% of the mixture, hence very negligible). For the mixtures at 200,000,
561 100,000, and 50,000 reads of BL21-DE3, the best metric was the mean abundance computed from the
562 node abundances without the abundances at zero, and by taking into account the multiple mapped reads.
563 While for the others, the best metric was the mean abundance computed from the node abundances
564 (including the abundances at zero), and by taking into account the multiple mapped reads.

565 This approach using linear models was particularly appropriate for this situation where the reference
566 variation graph and the sample contained a small number of strains and thus a small number of predictors
567 for the model. However, this can hardly transpose to a whole metagenomic sample with various species
568 and various strains that would lead to too many predictors and probably confusing the heuristics behind
569 the models. This was confirmed by applying the same methodology above on the mock dataset leading
570 to abundances estimation hardly comparable to expected. Compared to `Kraken2` results, the sum of
571 squared errors of our methodology was approximately 6 whereas for the results with the LASSO model it
572 was around 236. Nevertheless, those results highlighted the relevance of (i) using a metric taking into
573 account the multiple mapped reads and not only the unique mapped reads, and (ii) using our metric of
574 abundance based on the node abundances over raw read counts.

575 **S1.4 Performances**

576 Our benchmarks were performed on the GenOuest platform on a machine with 48 Xeon E5-2670 2.30
 577 GHz with 500 GB of memory and 16 CPUs. Time results (Table S1) are the wall-clock times. We
 578 provided rough computation time, mainly in the purpose to show that *StrainFLAIR* can be applied on
 579 usual datasets.

Dataset	Step	Items processed	Time	Disk used (GB)	Max mem. (GB)
Simulated	Gene prediction	7 genomes	0m20	0	1.2
	Gene clustering	34,011 genes	0m22	0	0.36
	Graph construction	8,596 clusters	2m44	0.04	1.31
	Graph concatenation	8,596 graphs	0m51	0	0.25
	Graph indexation	1 graph	6m23	0.16	4.24
	Mapping reads	350,000 short reads	15m15	0.16	0.99
	JSON conversion	1 GAMP file	3m58	4.2	0.03
	JSON parsing	1 JSON file + 1 GFA file + 1 pickle file	12m44	0	0.55
	Abundance computing	1 Gene abundances table	0m2	0	0.04
Mock	Gene prediction	91 genomes	1m43	1.02	6.7
	Gene clustering	280,174 genes	3m38	0.14	0.98
	Graph construction	270,712 clusters	41m54	1.12	9.1
	Graph concatenation	270,712 graphs	14m38	0	1.05
	Graph indexation	1 graph	75m19	1.98	30.4
	Mapping reads	21,389,196 short read pairs	147m28	7	17.5
	JSON conversion	1 GAMP file	53m21	75	0.12
	JSON parsing	1 JSON file + 1 GFA file + 1 pickle file	110m44	0	5.7
	Abundance computing	1 Gene abundances table	0m4	0	0.68

Table S1. *StrainFLAIR* performances on simulated and mock datasets.

580 **S1.5 Distance between the selected genomes in the simulated experiment**

581 We estimated the distance between the complete genomes of the selected strains using fastANI (Average
582 Nucleotide Identity). FastANI uses an alignment-free algorithm to estimate the average nucleotide identity
583 between pairs of sequences.

	K-12	IAI39	O104:H4	Sakai	SE15	Santai	BL21-DE3	RM8426
K-12	100	97.0652	98.3769	97.8703	96.8716	98.0362	98.9365	98.3657
IAI39	97.037	100	96.9742	96.7417	97.1289	96.9295	97.0197	96.8987
O104:H4	98.3059	96.9521	100	97.4788	96.8007	97.8896	98.249	98.7212
Sakai	97.7497	96.8627	97.5094	100	96.6657	98.1523	97.7455	97.6125
SE15	96.8453	97.1064	96.9211	96.7362	100	96.7575	96.8141	96.7763
Santai	98.0073	97.0372	97.9584	98.1797	96.8199	100	97.9279	97.9077
BL21-DE3	98.9983	97.1721	98.4048	97.8227	96.8448	97.9616	100	98.3204
RM8426	98.306	96.9037	98.6801	97.5815	96.6907	97.8353	98.2567	100

Table S2. Distance between each pair of complete genome sequences from eight strains of *E. coli* as computed by fastANI.

584 All pairs showed a distance at least greater than 95%, highlighting the strong similarities between
585 the strains. As a threshold, we although considered that beyond 99%, sequences were too similar to be
586 considered and distinguished, additionally to the effect of sequencing errors. The fastANI results showed
587 that none of the pairs exceeded this similarity threshold.

588 The strain *E. coli* BL21-DE3 was chosen as the unknown strain while the seven others would be used
589 to build the reference pangenome graph. According to the results of fastANI, the strain BL21-DE3 closest
590 genome in the present references is the strain K-12 with a similarity of 98.9%. Hence we expected to find
591 evidences of the strain K-12 while analyzing a sample containing the unknown strain BL21-DE3.

592 **S1.6 Detailed results from simulated datasets**

#reads K-12	Method	O104:H4	IAI39	K-12	Sakai	SE15	Santai	RM8426
1,000	Expected	59.88	39.92	0.2	0	0	0	0
	StrainFLAIR	56.45	43.55	0	0	0	0	0
	Kraken2	38.91	60.72	0.22	0.04	0.07	0.03	0.02
5,000	Expected	59.41	39.6	0.99	0	0	0	0
	StrainFLAIR	54.89	42.46	2.65	0	0	0	0
	Kraken2	38.61	60.25	0.99	0.04	0.07	0.03	0.02
10,000	Expected	58.82	39.22	1.96	0	0	0	0
	StrainFLAIR	54.08	41.96	3.96	0	0	0	0
	Kraken2	38.26	59.69	1.9	0.04	0.07	0.03	0.02
25,000	Expected	57.14	38.1	4.76	0	0	0	0
	StrainFLAIR	52.1	40.58	7.32	0	0	0	0
	Kraken2	37.23	58.1	4.51	0.04	0.07	0.03	0.02
50,000	Expected	54.55	36.36	9.09	0	0	0	0
	StrainFLAIR	49.23	38.51	12.26	0	0	0	0
	Kraken2	35.63	55.6	8.62	0.04	0.07	0.03	0.02
100,000	Expected	50	33.33	16.67	0	0	0	0
	StrainFLAIR	44.66	35.05	20.29	0	0	0	0
	Kraken2	32.8	51.19	15.85	0.04	0.07	0.03	0.02
200,000	Expected	42.86	28.57	28.57	0	0	0	0
	StrainFLAIR	38.12	29.83	32.05	0	0	0	0
	Kraken2	28.31	44.18	27.35	0.04	0.08	0.03	0.02

Table S3. Reference strains relative abundances expected and computed by StrainFLAIR or Kraken2 for each simulated experiment with variable coverage of the K-12 MG1655 strain. Best results are shown in bold.

593 Table S3 provides exhaustive results on simulated datasets when all queried strains are indexed in the
 594 variation graph. Table S4 provides exhaustive results on simulated datasets when one of the queried strain
 595 (BL21-DE3) is not indexed and highly similar to strain K-12.

#reads BL21-DE3	Method	O104:H4	IAI39	K-12	Sakai	SE15	Santai	RM8426
1,000	Expected	59.88	39.92	(0.2)	0	0	0	0
	StrainFLAIR	56.47	43.53	0	0	0	0	0
	Kraken2	38.93	60.76	0.11	0.05	0.08	0.04	0.03
5,000	Expected	59.41	39.6	(0.99)	0	0	0	0
	StrainFLAIR	56.45	43.55	0	0	0	0	0
	Kraken2	38.72	60.42	0.5	0.09	0.13	0.08	0.07
10,000	Expected	58.82	39.22	(1.96)	0	0	0	0
	StrainFLAIR	56.45	43.55	0	0	0	0	0
	Kraken2	38.47	60.05	0.92	0.14	0.19	0.12	0.13
25,000	Expected	57.14	38.1	(4.76)	0	0	0	0
	StrainFLAIR	54.09	41.71	4.2	0	0	0	0
	Kraken2	37.75	58.93	2.16	0.28	0.34	0.25	0.29
50,000	Expected	54.55	36.36	(9.09)	0	0	0	0
	StrainFLAIR	52.74	40.62	6.65	0	0	0	0
	Kraken2	36.59	57.17	4.15	0.51	0.57	0.48	0.53
100,000	Expected	50	33.33	(16.67)	0	0	0	0
	StrainFLAIR	50.47	38.64	10.89	0	0	0	0
	Kraken2	34.53	54.03	7.68	0.91	0.98	0.91	0.96
200,000	Expected	42.86	28.57	(28.57)	0	0	0	0
	StrainFLAIR	46.95	35.34	17.72	0	0	0	0
	Kraken2	31.14	48.83	13.53	1.57	1.67	1.58	1.68

Table S4. Reference strains relative abundances expected and computed by StrainFLAIR or Kraken2 for each simulated experiment with variable coverage of the BL21-DE3 strain, absent from the reference graph. BL21-DE3 being similar at 98.9% to K-12 strain (highest similarity compared to the other references), we expect that reads from BL21-DE3 will map this strain, hence its expected values are given in parentheses, as they correspond to BL21-DE3 strain abundances and not K-12. Best results are shown in bold.

