

1 **MetaPGN: a pipeline for construction and graphical visualization of** 2 **annotated pangenome networks**

3 4 **Abstract**

5 Pangenome analyses facilitate the interpretation of genetic diversity and evolutionary history of a
6 taxon. However, there is an urgent and unmet need to develop new tools for advanced pangenome
7 construction and visualization, especially for metagenomic data. Here we present an integrated
8 pipeline, named MetaPGN, for construction and graphical visualization of pangenome network
9 from either microbial genomes or metagenomes. Given either isolated genomes or metagenomic
10 assemblies coupled with a reference genome of the targeted taxon, MetaPGN generates a
11 pangenome in a topological network, consisting of genes (nodes) and gene-gene genomic
12 adjacencies (edges) of which biological information can be easily updated and retrieved. MetaPGN
13 also includes a self-developed Cytoscape plugin for layout of and interaction with the resulting
14 pangenome network, providing an intuitive and interactive interface for full exploration of genetic
15 diversity. We demonstrate the utility of MetaPGN by constructing *Escherichia coli* (*E. coli*)
16 pangenome networks from five *E. coli* pathogenic strains and 760 human gut microbiomes
17 respectively, revealing extensive genetic diversity of *E. coli* within both isolates and gut microbial
18 populations. With the ability to extract and visualize gene contents and gene-gene physical
19 adjacencies of a specific taxon from large-scale metagenomic data, MetaPGN provides advantages
20 in expanding pangenome analysis to uncultured microbial taxa. MetaPGN is available at
21 <https://github.com/peng-ye/MetaPGN>.

22 **Keywords:** pangenome, visualization, metagenomics

23

24

25

26

27 Introduction

28 The concept of the pangenome, defined as the full complement of genes in a clade, was first
29 introduced by Tettelin *et al.* in 2005 [1]. Pangenome analyses of a species now provide insights
30 into core- and accessory-genome profiles, within-species genetic diversity, evolutionary dynamics
31 and niche-specific adaptations. A number of methods and tools have to date been proposed for
32 pangenome analysis on genomic or metagenomic data (Table 1).

33 Typical pangenome tools such as GET_HOMOLOGUES [2] and PGAP [3], mainly focus on
34 analyzing homologous gene families and calculating the core/accessory genes of a given taxon.
35 However, these tools cannot provide the variations of gene-gene physical relationships. Tools like
36 GenoSets [4], PGAT [5], PEGR [6], EDGAR [7], GenomeRing [8] and PanViz [9] are developed
37 to generate a linear or circular presentation of compared genomes, which can indicate the physical
38 relationships between genomic sequences or genes. However, in the linear or circular
39 representations generated by these tools, the same homologous region is visualized multiple times
40 and shown on separate input genomes. Hence, it will be difficult for users to track a homologous
41 region among the input genomes, especially when there is a large number of homologous regions
42 and input genomes.

43 Pangenomes built using *de Bruijn* graph, like SplitMEM [10] and a tool introduced by Baier *et*
44 *al.* [11], partly solve the above problems. In the resulting graph generated by these tools, the
45 complete pangenome is represented in a compact graphical representation such that the
46 core/accessory status of any genomic sequences is immediately identifiable, along with the context
47 of the flanking sequences. This strategy enables powerful topological analysis of the pangenome
48 not possible from a linear/circular representation. Nevertheless, tools based on the *de Bruijn* graph
49 algorithm can only construct a compact network comprised of core/accessory genomic sequences
50 instead of genes, which means retrieving or updating functional information in downstream
51 analysis will be difficult. Furthermore, these tools do not visualize the constructed *de Bruijn* graph
52 and provide an interactive interface for users to explore the graph.

53 Moreover, all the above-mentioned tools analyze pangenomes via genomic data which require
54 organisms isolated from the environment and cultured *in vitro*. Recent advances in metagenomics
55 have led to a paradigm shift in pangenome studies from a limited quantity of cultured microbial

56 genomes to large-scale metagenomic datasets containing huge potential for functional and
57 phylogenetic resolution from the still uncultured taxa. Several existing tools dealing with
58 metagenomic data are based on constructed pangenomes and cannot utilize the abundant gene
59 resources contained in metagenomes to extend the pangenomes in question. For example,
60 PanPhlAn [12], MIDAS [13], and a pipeline introduced by Delmont and Eren [14] maps reads
61 onto a reference pangenome, to describe the pattern of the presence/absence of genes in
62 metagenomes. As for another example, Kim *et al.* [15] clustered genes predicted from
63 metagenomic contigs with *Bacillus* core genes for profiling the *Bacillus* species in the
64 microbiomes. Recently, Farag *et al.* [16] aligned metagenome contigs with reference genomes for
65 identification of “*Latescibacteria*” genomic fragments. Even though this strategy can theoretically
66 recruit sequences not present in the reference genomes, it is likely to filter out “*Latescibacteria*”
67 genomic fragments with structural variations compared to the reference ones. Furthermore, all
68 these aforementioned methods using metagenomic data do not organize the pangenome using a
69 network, which is essential for efficiently storage and visualization of pangenomes constructed
70 from metagenomic data.

71 Here, we introduce an integrated pipeline (MetaPGN) for network-based construction and
72 visualization of prokaryotic pangenomes for both isolated genomes and metagenomes. Given
73 genomic or metagenomic assemblies and a reference genome of a taxon of interest, MetaPGN
74 derives a pangenome network for integrating genes (nodes) and gene-gene adjacencies (edges)
75 belonging to a given taxon. MetaPGN also includes a specific Cytoscape plugin for layout of and
76 interaction with the resulting pangenome network, providing an intuitive and interactive interface
77 for the exploration of gene diversity. For example, in the visualized network in Cytoscape, users
78 can specify gene annotations, customize the appearance of nodes and edges, and search and
79 concentrate on genes of certain functions. We applied MetaPGN on assemblies from five
80 pathogenic *E. coli* strains and 760 human gut microbiomes respectively, with *E. coli* K-12 substr.
81 MG1655 (*E. coli* K-12) being the reference genome. Our results showed that by taking gene
82 adjacency into account and visualizing the pangenome network in a well-organized manner,
83 MetaPGN can assist in illustrating genetic diversity in genomic or metagenomic assemblies
84 graphically and conveniently.

85

Table 1. Comparison of several pangenome analysis methods.

Method	Input		Output			Functionality		
	Isolate genomes	Metagenomes	Gene content	Gene-gene adjacency	Network	Biological annotation	Interactive visualization	
GET_HOMOLOGUES [2] and PGAP [3]	Yes	No	Yes	No	No	Yes	No	
GenoSets [4], PGAT [5], PEGR [6], EDGAR [7], GenomeRing [8]	Yes	No	Yes	Yes	No	Yes	No	
PanViz [9]	Yes	No	Yes	Yes	No	Yes	Yes	
SplitMEM [10] and a tool introduced by Baier <i>et al.</i> [11]	Yes	No	Yes	Yes	Yes	No	Yes	
PanPhlAn [12], MIDAS [13] and a method introduced by Farag <i>et al.</i> [16]	No	Yes	Yes	No	No	Yes	No	
MetaPGN	Yes	Yes	Yes	Yes	Yes	Yes	Yes	

87 Results

88 **General workflow.** MetaPGN accepts genome or metagenome assemblies as input (query
89 assemblies) and requires a reference genome for recruitment of the query assemblies and as the
90 skeleton of the pangenome network. The MetaPGN pipeline can be divided into two main parts:
91 (i) construction of a pangenome network comprised of representative genes, including gene
92 prediction, gene redundancy elimination, gene type determination, pairwise gene adjacency
93 extraction, assembly recruitment (for metagenomic assemblies), and pangenome network
94 generation, and (ii) visualization of the pangenome network in an organized way, where nodes
95 represent genes and edges indicate gene adjacencies, in Cytoscape [17] with a self-developed
96 plugin (Fig. 1, Fig. S1, and Methods). From the resultant pangenome network, the degree of
97 similarity among homologous genes, as well as their genomic context is easily visible. Of note,
98 users can further add and update annotation for nodes and edges in the networks, based on which
99 elements of interest can be accessed conveniently.

100

101 **Pangenome network of 5 pathogenic *Escherichia coli* genomes.** In order to demonstrate its
102 potential in studying microbial genetic diversity and phenotype-genotype relationship, we first
103 applied MetaPGN on genomes of 5 pathogenic *E. coli* isolates, *E. coli* O26:H11 str. 11368, *E.*
104 *coli* O127:H6 E2348/69, *E. coli* O157:H7 str. EDL933, *E. coli* O104:H4 str. 2011C-3493 and *E.*
105 *coli* 55989. A commensal *E. coli* strain, K-12 substr. MG1655 (Supplementary Table S1) was
106 chosen as the reference genome in this instance and in all examples shown below.

107 A pangenome network consisting of 9,161 nodes and 11,788 edges (Supplementary Table S3,
108 Supplementary File 2) was constructed and visualized (Methods). Based on the well visualized
109 pangenome network along with functional annotation, we can now graphically observe the extent
110 of variations of certain genes, as well as their genomic context. For example, when focusing on a
111 cluster of flagellar genes (Fig. 2a), we found that *fliC* sequences encoding the filament structural
112 protein (H-antigen) and *fliD* sequences encoding the filament capping protein are highly divergent
113 with nucleotide sequence identity < 95% and/or overlap < 90% among these *E. coli* strains (See
114 Methods). In contrast, four genes encoding chaperones (*fliS*, *fliT*, *fliY*, *fliZ*) and a gene related to
115 regulation of expression of flagellar components (*fliA*) are conserved (nucleotide sequence identity

116 $\geq 95\%$ and overlap $\geq 90\%$) over all the *E. coli* strains investigated. A gene (270bp) encoding a
117 hypothetical protein is uniquely presented between *fliC* and *fliD* in *E. coli* O157:H7 str. EDL933.

118 In a fimbria protein-related gene cluster, compared to the reference *E. coli* strain, all the 5
119 pathogenic strains possess several genes located between two conserved genes encoding an outer
120 membrane protein and a regulatory protein, and *E. coli* O127:H6 E2348/69 uniquely exhibits more
121 genes encoding proteins of unknown functions (Fig. 2b).

122 For a gene cluster responsible for the biosynthesis of lipopolysaccharides (LPS), *E. coli*
123 O127:H6 E2348/69 shares three genes with the reference strain that differentiate from the other 4
124 pathogenic strains (Fig. 2c). For another gene cluster of related function, the *E. coli* O127:H6
125 E2348/69 also shows a strain-specific duplication event of two genes involved in colanic acid (CA)
126 synthesis (*wcaH* and *wcaG*, denoted by a purple dash line in Fig. 2d). It has been demonstrated
127 that CA can modify lipopolysaccharide (LPS) generating a novel form (M_{LPS}) which may enhance
128 survival of *E. coli* in different ways [18]. The two *wcaH* genes in *E. coli* O127:H6 E2348/69, may
129 even though they share high similarity (99.1% identity) confer the strain with different functional
130 potentials for CA formation and thereby novel survival mechanisms.

131 In addition, the German outbreak *E. coli* O104:H4 str. 2011C-3493 shares identical nodes and
132 edges in the flagellar-related gene cluster (Fig. 2a) and the O antigen-related gene cluster with a
133 historical *E. coli* 55989 (Fig. 2d), suggesting a close
134 evolutionary relationship between these strains as previously reported [19,20].

135 These results demonstrate the feasibility of MetaPGN for construction and visualization of
136 microbial pangenomes in an organized way. Moreover, by involving genomic adjacency and
137 offering easy-to-achieve biological information, MetaPGN provides a convenient way to assist
138 biologists in exposing genetic diversity for genes of interest among the organisms under study.

139

140 **Pangenome network of *E. coli* in 760 metagenomes.** Moving beyond surveying the pangenome
141 network of isolate genomes, we applied MetaPGN in metagenomic datasets to interrogate the *E.*
142 *coli* pangenome network on a grander scale. Assemblies of 760 metagenomes sequenced in the
143 Metagenomics of the Human Intestinal Tract (MetaHIT) project [21–24] were collected, which
144 contained 8,096,991 non-redundant genes with annotations [24]. As metagenome assemblies are
145 from varied taxa, it is necessary to recruit assemblies of the targeted taxon before construction of
146 the pangenome network. In this study, metagenome assemblies were recruited using a gene

147 alignment-based strategy, which was assessed with mock datasets (Methods). With the recruited
148 assemblies, a pangenome network consisting of 9,406 nodes and 14,676 edges (Supplementary
149 Table S3, Supplementary File S3) was generated and visualized after refinement (Methods).

150 Based on annotation, we first searched flagellin-related genes in this network. We found that
151 the pattern of adjacencies among these genes was similar to that in the pangenome network of the
152 5 pathogenic *E. coli* genomes: *fliC* and *fliD* are hypervariable while *fliT*, *fliY*, *fliZ* and *fliA* are very
153 conserved among these 760 samples. However, some genes of unknown function locate between
154 *fliC* and *fliA* (Fig. 3a), instead of between *fliC* and *fliD* in the pangenome network of the 5
155 pathogenic *E. coli* strains (Fig. 2a).

156 We then investigated mobile genetic elements (MGEs) in this pangenome network, as they can
157 induce various types of genomic rearrangements[25]. Of the 362 nodes (~4%) annotated as MGE-
158 related (according to Cluster of Orthologous Groups annotation done in reference [24]), many were
159 flanked by shared genes on different *E. coli* genomes. In a region of the network, a gene cluster
160 containing MGEs is query-specific, indicating there might be genomic rearrangements caused by
161 strain-specific MGEs within the *E. coli* species (Fig. 3b). In another part of the network harboring
162 MGEs, we observed that several branches of non-MGE genes are inserted in between two MGEs,
163 which may imply a mutation hotspot within the region, or the existence of MGEs as yet
164 undescribed (Fig. S1).

165 Application of MetaPGN in large-scale metagenomic data generated an *E. coli* pangenome
166 network that might hardly be constructed from isolated genomes. As demonstrated here, the
167 assembly-recruitment based, well-organized and visualized pangenome network can greatly
168 expand our understanding in the genetic diversity of a taxon, although further efforts in
169 bioinformatic and experimental analyses are needed to verify and extend these findings.

170

171 **Assessment of pangenome networks derived from metagenomes.** Affected by the complexity
172 of microbial communities, limitations in sequencing platforms and imperfections of bioinformatic
173 algorithms, a genomic sequence of an organism is frequently split into dozens of assemblies when
174 assembled from metagenomic reads. Due to this nature, a pangenome network recovered from a
175 limited number of assemblies is likely to be segmented compared to a complete genome. To
176 propose a minimum size of assemblies for getting an approximately complete connected

177 pangenome network, we assessed the completeness of *E. coli* pangenome networks derived from
178 varying size of recruited assemblies (Methods). As shown in Fig. 4, the count of connected
179 subnetworks drops dramatically with the total length of recruited assemblies increasing from 5 Mb
180 to 50 Mb (roughly from 1 × to 10 × of a *E. coli* genome), then barely changes even when using
181 all recruited assemblies of the dataset (215 Mb, from 760 samples). Based on this analysis, a
182 minimum size of recruited assemblies 10-fold of the studied genome is required to generate a
183 relatively intact pangenome network when constructed from metagenomes.

184

185 Discussion

186 Since first coined more than a decade ago, pangenome analysis has provided a framework for
187 studying the genomic diversity within a species. Current methods for pangenome analyses mainly
188 focus on gene contents but ignore their genomic context, as well as having shortages in pangenome
189 visualization. Besides, available methods are usually designed for genomic data and not capable
190 of constructing pangenomes from metagenomics data. To fill these gaps, our MetaPGN pipeline
191 takes genome or metagenome assemblies as input, uses gene contents as well as pairwise gene
192 adjacency to generate a compact graphical representation for the gene network based on a reference
193 genome, and visualizes the network in Cytoscape with a self-developed plugin (Fig. 1, Fig. S2).

194 From the two MetaPGN-derived *E. coli* pangenome networks, we can directly observe the
195 diversity of genes among the five pathogenic *E. coli* strains and 760 human gut microbiomes with
196 respect to the reference genome. For instance, in the pangenome network for the 5 pathogenic *E.*
197 *coli* strains, we found that nucleotide sequences of the *fliC* gene which carries H-antigen specificity
198 were highly divergent among the *E. coli* assemblies (Fig. 2a). These *fliC* sequences were more
199 varied in the 760 human gut microbiomes (Fig. 3a). In addition, genes for synthesis of O-antigen
200 and outer membrane protein showed a great diversity in the pangenome network of the 5 *E. coli*
201 strains (Fig. 2c, Fig. 2d). These results are in agreement with previous findings on H-antigen
202 specificity related genes [26–28] and O-antigen related genes [29,30]. We also showed that when
203 gene adjacency is incorporated into the construction and visualization of pangenomes, locations
204 of genes of unknown function are identified, which may be helpful for the inference of their
205 biological functions. For example, in both the two pangenome networks, we found genes of
206 unknown function locating between the *fliC* gene and other flagellin-related genes (Fig. 2a, located

207 between *fliC* and *fliD*, Fig. 3a, and located between *fliC* and *fliA*), indicating that these functionally
208 unknown genes may play a role in flagellin biosynthesis [31], although further experimental trials
209 are needed to prove this point. Additionally, from the pangenome network of the five *E. coli* strains,
210 we observed a variation in *E. coli* O127:H6 E2348/69, which was shown to stem from a duplication
211 event of two genes involved in colanic acid synthesis (*wcaH* and *wcaG*, Fig 2d). This finding
212 indicates that knowledge of genomic adjacency may also shed light on structural variations among
213 the input assemblies. If extended, genomic adjacency may further help in finding possible
214 functional sequences which are associated with structural variations, as Delihias [32] and Wang *et*
215 *al.* [33] reported on repeat sequences concentrated at the breakpoints of structural variations.
216 Studying genomic adjacency can also improve the discovery of potential functional modules, as
217 Doron *et al.* [34] systematically discovered bacterial defensive systems by examining gene
218 families enriched next to known defense genes in prokaryotic genomes. These examples illustrate
219 the value of including gene adjacencies in visualizing a pangenome to retrieve biological
220 information. Although the examples shown in this study use the genome of a commensal *E. coli*
221 strain for assembly recruitment and network arrangement, users can specify the reference genome
222 when applying MetaPGN. Epidemiologists can use MetaPGN to compare assemblies of outbreak
223 strains or viruses, such as *Vibrio cholerae* or Ebola virus, with those of some well-studied
224 pathogenic strains to find novel variations involved in pathogenesis, which may further provide
225 candidate targets for drug and vaccine design [35,36].

226 Genomic variants of intestinal bacteria were found to be correlated with diseases. As one
227 example, among the common members of the normal colonic microbiota, *Bacteroides fragilis* (*B.*
228 *fragilis*), the inclusion of a pathogenicity island (BfPAI) distinguished enterotoxigenic strains
229 (ETBF) from nontoxigenic ones (NTBF), by their ability to secrete a zinc-dependent
230 metalloprotease toxin that can induce inflammatory diarrhea and even colon carcinogenesis
231 [37,38]. As another example, Scher *et al.* performed shotgun sequencing on fecal samples from
232 newly-onset untreated rheumatoid arthritis (NORA) patients and healthy individuals, and
233 identified several NORA-specific *Prevotella copri* genes [39]. Hence, pangenome networks built
234 from metagenomes of patients and healthy subjects may aid in detecting associated or causal
235 genomic variants of a certain species.

236 It should be noticed that, in this pipeline, we compare genes depending on nucleotide-level
237 sequence identity and overlap: genes with $\geq 95\%$ identity and $\geq 90\%$ overlap are regarded as the
238 same gene. However, genes sharing the same function may not satisfy this criterion ($\geq 95\%$ identity
239 and $\geq 90\%$ overlap), and protein encoded by these genes may exhibit more similarity due to
240 different codon usage. Hence, in our future work, we intend to cluster genes by comparing their
241 nucleotide sequences as well as the amino acid sequences. Furthermore, the current MetaPGN
242 pipeline does not consider other genomic features or physical distances between genes in
243 constructing the pangenome network. Thus, differences in other genomic features such as
244 ribosomal binding site (RBS) sequences [40,41] and distances between the RBS and start codons
245 [42] may result in distinct phenotypes. Accordingly, users may include such information in
246 analyzing pangenome networks.

247 To conclude, MetaPGN enables direct illustration of genetic diversity of a species in pangenome
248 networks, improving understanding of genotype-phenotype relationships and evolutionary history.
249

250 Methods

251 **Pangenome network construction in MetaPGN.** First, gene prediction of query assemblies is
252 performed using MetaGeneMark (Version 2.8) [43]. In order to eliminate redundancy, the resultant
253 genes are clustered by CD-HIT (Version 4.5.7) [44] with identity $\geq 95\%$ and overlap ≥ 90 , and
254 genes in a same cluster are represented by the longest sequence of the cluster which is termed the
255 representative gene. Representative genes of all clusters are subsequently aligned against genes on
256 the given reference genome using BLAT (Version 34) [45]. From the alignment result, genes
257 shared between the representative gene set and the reference gene set with identity $\geq 95\%$
258 and overlap $\geq 90\%$ are defined as ‘shared genes’. The remaining representative and reference genes
259 other than those shared genes are defined as ‘query-specific genes’ and ‘reference-specific genes’,
260 respectively. Pairwise gene physical adjacency of representative genes on the query assemblies
261 and of reference genes are then extracted, and status for each adjacency of being ‘shared’, ‘query-
262 specific’, or ‘reference-specific’ is determined. Finally, based on the recruited assemblies and the
263 reference genome, an initial pangenome network is generated: each node stands for a reference
264 gene or a representative gene on the recruited assemblies; two nodes are connected by an edge if
265 they are physically adjacent on the recruited assemblies or on the reference genome. The weight

266 of a node or an edge denotes its occurrence frequency on all of the recruited assemblies and the
267 reference genome.

268

269 **Pangenome network visualization in MetaPGN.** The following preprocessing work on the
270 initial pangenome network was implemented before visualization: 1. The initial pangenome
271 network was refined by removing isolated networks (networks not connected with the backbone)
272 and tips (nodes only connected with another node); 2. Nodes and edges were added with some
273 extra attributes, such as the status of the nodes and edges (query-specific, reference-specific or
274 shared), whether the genes for the nodes were phage-, plasmid-, CRISPR- related genes and so on
275 (Supplementary Table S3). Users can specify the attributes of nodes and edges according to their
276 own datasets.

277 We then used a self-developed Cytoscape plugin to visualize the pangenome network in an
278 organized way (Supplementary Text 2 in Supplementary File S1 illustrates how to install and use
279 the plugin in Cytoscape). Our algorithm for organizing nodes in the network is as follows:

- 280 1. Construct a circular skeleton for the pangenome network with shared nodes and reference-
281 specific nodes, according to positions of their related reference genes on the reference genome.
282 If there are two or more representative genes similar to the same reference gene ($\geq 95\%$ identity
283 and $\geq 90\%$ overlap), use one of these representative genes to construct the skeleton and place
284 the others on both sides of the skeleton in turn (Fig. S2 a).
- 285 2. Arrange query-specific nodes region by region, including,
 - 286 2.1. Select query-specific nodes in a region spanning less than 30 nodes in the skeleton
287 (see Supplementary Text 3 in Supplementary File S2 for more details).
 - 288 2.2. Arrange these query-specific nodes as follows,
 - 289 i. For those that directly link with two nodes on the skeleton, place them on the bisector
290 of the two skeleton nodes. If there are two or more query-specific nodes directly
291 linking with the same pair of nodes on the skeleton, place them on both sides of the
292 bisector of these pair of skeleton nodes in turn (Fig. S2 b).
 - 293 ii. Among the remaining nodes, for those that directly link with two placed nodes, place
294 them on the bisectors of the placed ones. Iterate this step for five times (Fig. S2 c).

295 iii. For the remaining nodes, place them into an arc without moving the placed nodes (Fig.
296 S2 d), or else place them one by one starting near a placed node (Fig. S2 e).

297

298 **Construction and visualization of the 5-*E. coli*-genome pangenome network.** Genes were
299 extracted from the complete genome for each strain (Supplementary Table S1). With *E. coli* K-12
300 as the reference, a pangenome network was generated for these five *E. coli* strains using our
301 MetaPGN tool. In the visualization of this pangenome network, we used green, blue and red color
302 to denote a reference-specific, shared, and query-specific node or edge, respectively, and specified
303 sizes of nodes and widths of edges with their occurrence frequency in the input genomes.

304

305 **Assessment of the gene alignment-based assembly recruitment strategy.** A gene alignment-
306 based strategy was used for recruitment of metagenome assemblies in this study, which considers
307 1) the count of genes on an assembly (c), and 2) the ratio of the number of shared genes (designated
308 as aforementioned) on an assembly to the total number of genes on that assembly (r). $c = 3$ paired
309 with $r = 0.5$, requiring at least 3 genes including 2 shared genes containing in an assembly, was
310 chosen for recruitment of metagenome assemblies in this study.

311 5 mock metagenomic datasets were used to assess the performance of this strategy. Briefly,
312 simulated reads of 60 bacterial genomes from 14 common genera (*Bifidobacterium*, *Clostridium*,
313 *Enterobacter*, *Escherichia*, *Haemophilus*, *Klebsiella*, *Lactobacillus*, *Neisseria*, *Pseudomonas*,
314 *Salmonella*, *Shigella*, *Staphylococcus*, *Streptococcus*, *Yersinia*) present in the human gut
315 (Supplementary Table S1), including the 5 pathogenic *E. coli* strains mentioned above and 10
316 strains from *E. coli*-closely-related species (*Enterobacter aerogenes*, *Enterobacter cloacae*,
317 *Escherichia albertii*, *Escherichia fergusonii*, *Klebsiella oxytoca*, *Klebsiella pneumoniae*, *Shigella*
318 *boydii*, *Shigella sonnei* and *Salmonella enterica*), were generated by iMESSi [46]. Each dataset
319 was simulated at the same complexity level with 100 million (M) 80bp paired-end reads of 12
320 strains from 11-12 different genera, including 2 strains of closely related species to *E. coli*, and the
321 relative abundances of strains were assigned by the broken-stick model (Supplementary Table S2).
322 Simulated reads were first independently assembled into assemblies by SOAPdenovo2 in each
323 dataset [43], with an empirical k-mer size of 41. Genes were then predicted on assemblies longer

324 than 500bp using MetaGeneMark [42] (default parameters were used except the minimum length
325 of genes was set as 100bp).

326 Assemblies of each mock dataset were first aligned against the 5 pathogenic *E. coli* reference
327 genomes by BLAT [45]. Those assemblies that have an overall $\geq 90\%$ overlap and $\geq 95\%$ identity
328 with the reference genomes were considered as *E. coli* genome-derived (traditional genome
329 alignment-based strategy). Those *E. coli* genome-derived assemblies containing at least three
330 genes (i.e., containing at least two edges) were recruited for construction of a reference pangenome
331 network (RPGN). A query pangenome network (QPGN) was then generated from assemblies
332 selected by the gene alignment-based strategy with $c = 3$ and $r = 0.5$ as described above.

333 Accuracy of query assembly recruitment was assessed, in respect of conformity and divergence
334 between the RPGN with the QPGN (Supplementary Text 4 and 5 in Supplementary File S2). The
335 result showed that the QPGN recovered 84.3% of node and 84.7% of edge in the RPGN, while
336 falsely included 1.1% of node and 2.2% of edge, which demonstrated the high accuracy of the
337 gene alignment-based strategy for recruitment of metagenome assemblies.

338

339 **Construction and visualization of the 760-metagenome pangenome network.** Assemblies and
340 representative genes of the 760 metagenomes generated in Reference [24] were used here, since
341 they were produced using identical methods and parameter settings in this study. A pangenome
342 network was generated following steps described above, again using *E. coli* K-12 as the reference,
343 and $c = 3$, $r = 0.5$ for assembly recruitment. The resulting pangenome network was visualized in
344 the same way as visualizing the 5-*E. coli*-genome pangenome network.

345

346 **Analysis of subnetworks comprising a pangenome network.** 10-700 metagenomes were
347 randomly sampled from the above-mentioned 760 metagenomes. For each sub-dataset, a
348 pangenome network was constructed after assembly recruitment using *E. coli* K-12 as the
349 reference genome. For each pangenome network, reference-specific edges were removed before
350 counting the number of subnetworks. Only sub-datasets with a size of recruited assemblies greater
351 than 5 Mb were used to generate the scatterplot, in which a curve with 95% confidence intervals
352 was fitted by the 'loess' smoothing method in R [47].

353

354 **Computational resources and runtime**

355 Timings for major steps of the MetaPGN pipeline are shown below. Tests were run on a single
356 CPU of an Intel Core Processor (Broadwell) processor with 64 GB of RAM, without otherwise
357 specified. The timings were CPU time including parsing input and writing outputs (h for hours, m
358 for minutes, and s for seconds).

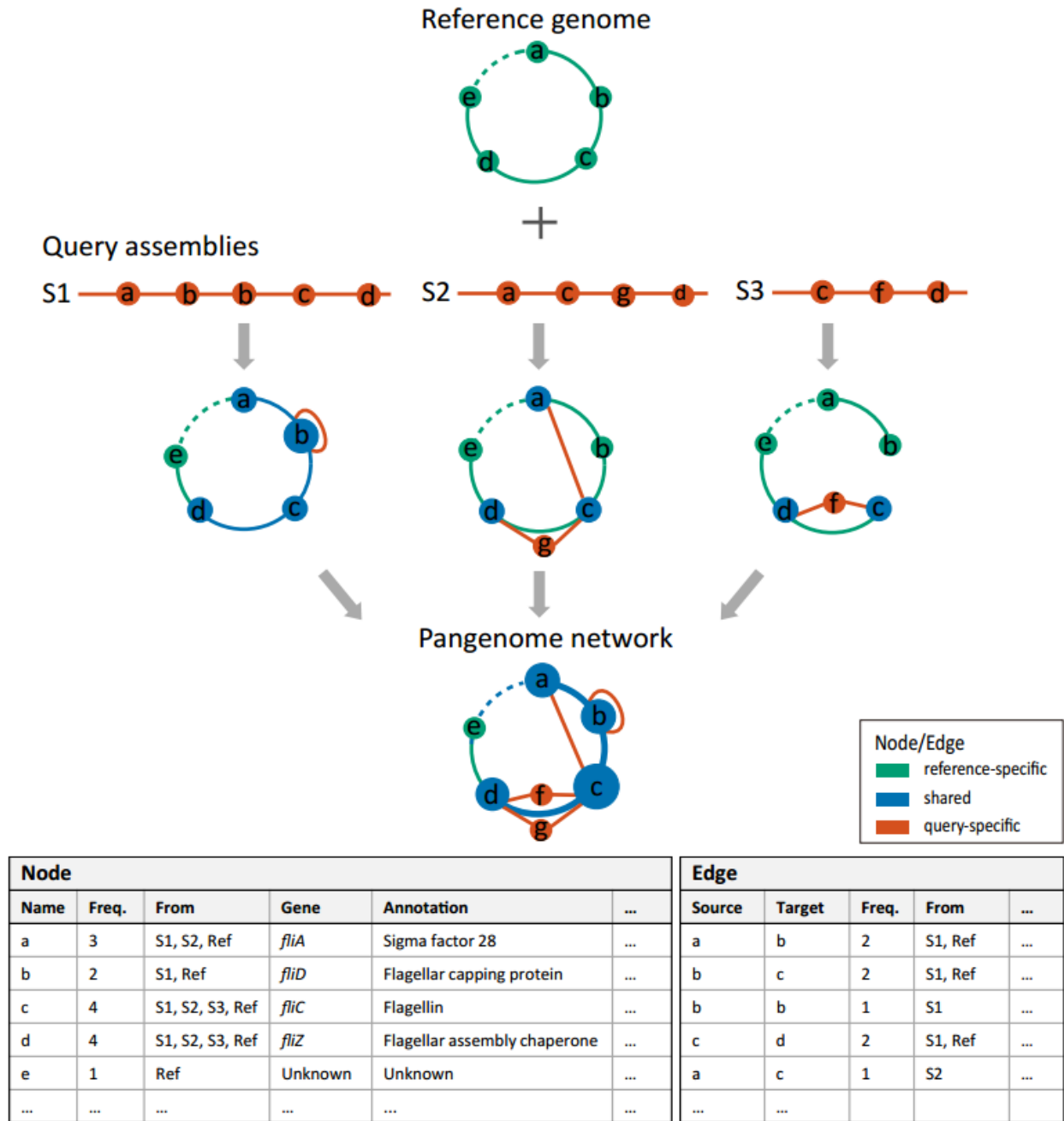
359 The average time for gene prediction for a mock metagenome was 7s, and it varies depending on
360 the size of a metagenome. The time for redundancy elimination of genes using CD-HIT [44] was
361 1m 44s for the 5 *E. coli* stains, 50m 19s for the 5 mock datasets. For the 760 metagenomes, to
362 perform redundancy elimination parallelly, we divided all genes into 200 sections, which resulted
363 in 20,101 [$N = (n + 1) \times (n \div 2) + 1, n = 200$] clustering tasks, and then submitted each task
364 onto available machines in a high-performance computing cluster. The dividing step took 20m 4s
365 with a peak memory usage of 10GB in the local machine, and the average time for a clustering
366 task was 44m with taking less than 3GB of RAM, consuming total time of 14,814h. The time for
367 recognizing the status (reference-specific, query-specific or shared) for nodes and edges was 10s
368 for the 5 *E. coli* strains, 1m for the 5 mock datasets and 24m for the 760 metagenomes. Finally,
369 the generation of the pangenome network took less than 1s for the 5 *E. coli* strains, less than 1s for
370 the 5 mock datasets and 3m 35s for the 760 metagenomes.

371

372 **Data availability.** Genome sequence of 60 strains (including 5 *E. coli* strains) and the *E. coli* K-
373 12 reference genome were downloaded from the National Center for Biotechnology Information
374 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/bacteria/>, Please refer to **Supplementary Table S1** for
375 detailed information). Sequencing data of the 760 metagenomes were previously generated in the
376 Metagenomics of the Human Intestinal Tract (MetaHIT) project [21–24], and assemblies of these
377 760 metagenomes are deposited at EBI under PRJEB28245. The MetaPGN pipeline, related
378 manuals and Cytoscape session files for *E. coli* pangenome networks derived from five pathogenic
379 *E. coli* strains and from 760 metagenomes are available on Github ([https://github.com/peng-](https://github.com/peng-ye/MetaPGN)
380 [ye/MetaPGN](https://github.com/peng-ye/MetaPGN)) and SciCrunch (SCR_016454).

381

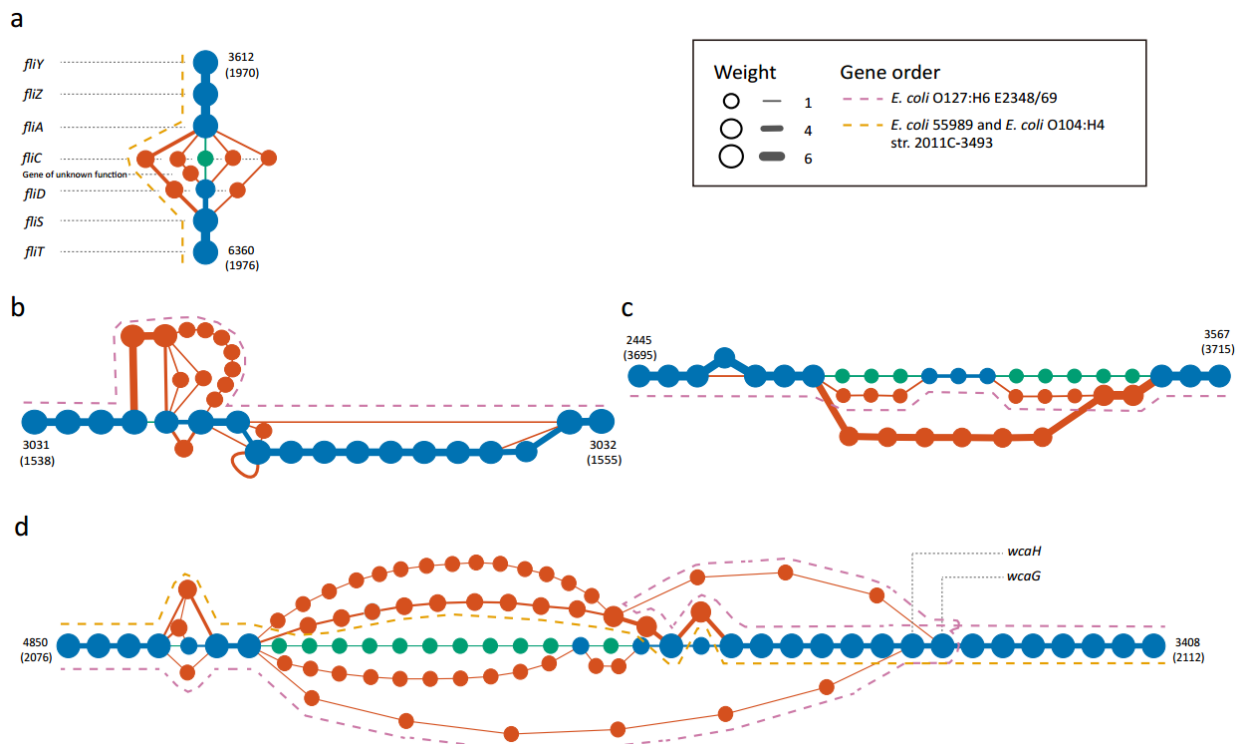
382 List of Figures



383 **Figure 1.** An Overview of the MetaPGN pipeline: from assemblies to a pangenome network. Gene
 384 prediction is performed on query assemblies. The resulting genes are clustered, after which genes
 385 in the same cluster are represented by the longest sequence of this cluster called the representative
 386 gene (node a-g). All these representative genes are then aligned against genes on the given
 387 reference genome. From the alignment result, genes shared between the representative gene set

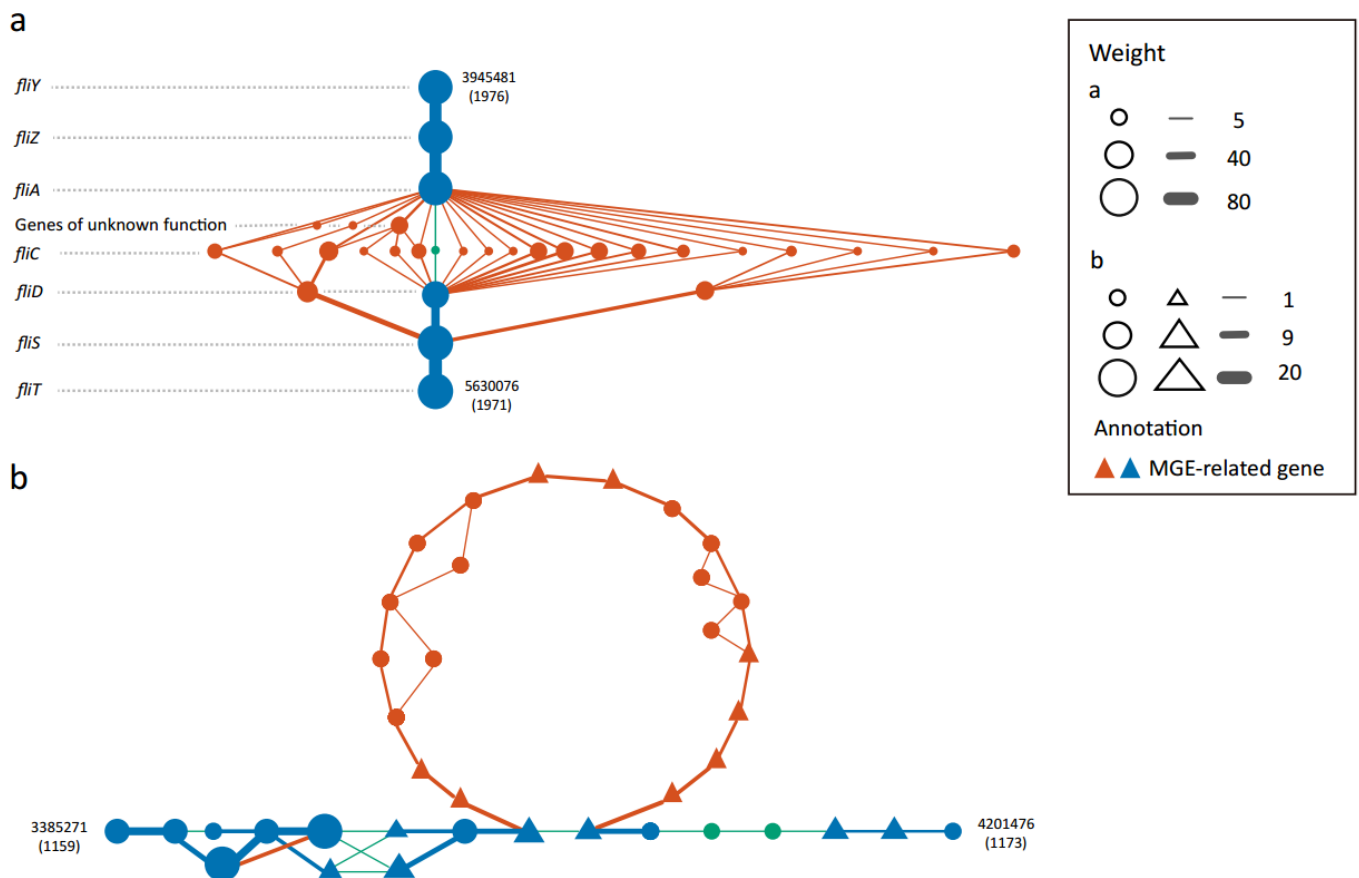
388 and the reference gene set are defined as ‘shared’ genes (blue). The remaining representative and
389 reference genes other than those shared genes are defined as ‘query-specific’ genes (red) and
390 ‘reference-specific’ genes (green), respectively. Pairwise gene physical adjacency of
391 representative genes on the query assemblies and of reference genes are then extracted, and status
392 for each adjacency of being ‘shared’ (blue), ‘query-specific’ (red), or ‘reference-specific’ (green)
393 is determined. Finally, based on the recruited assemblies and the reference genome, a pangenome
394 network is generated: each node stands for a reference gene or a representative gene on the
395 recruited assemblies; two nodes are connected by an edge if they are physically adjacent on the
396 recruited assemblies or the reference genome. The weight of a node or an edge is its occurrence
397 frequency on all of the recruited assemblies and the reference genome (Methods). The pangenome
398 network is then visualized in Cytoscape with a self-developed plugin (Methods) for a better
399 arrangement. Biological information of nodes and edges, such as gene name and annotation, can
400 be easily retrieved in the interactive user interface in Cytoscape.

401

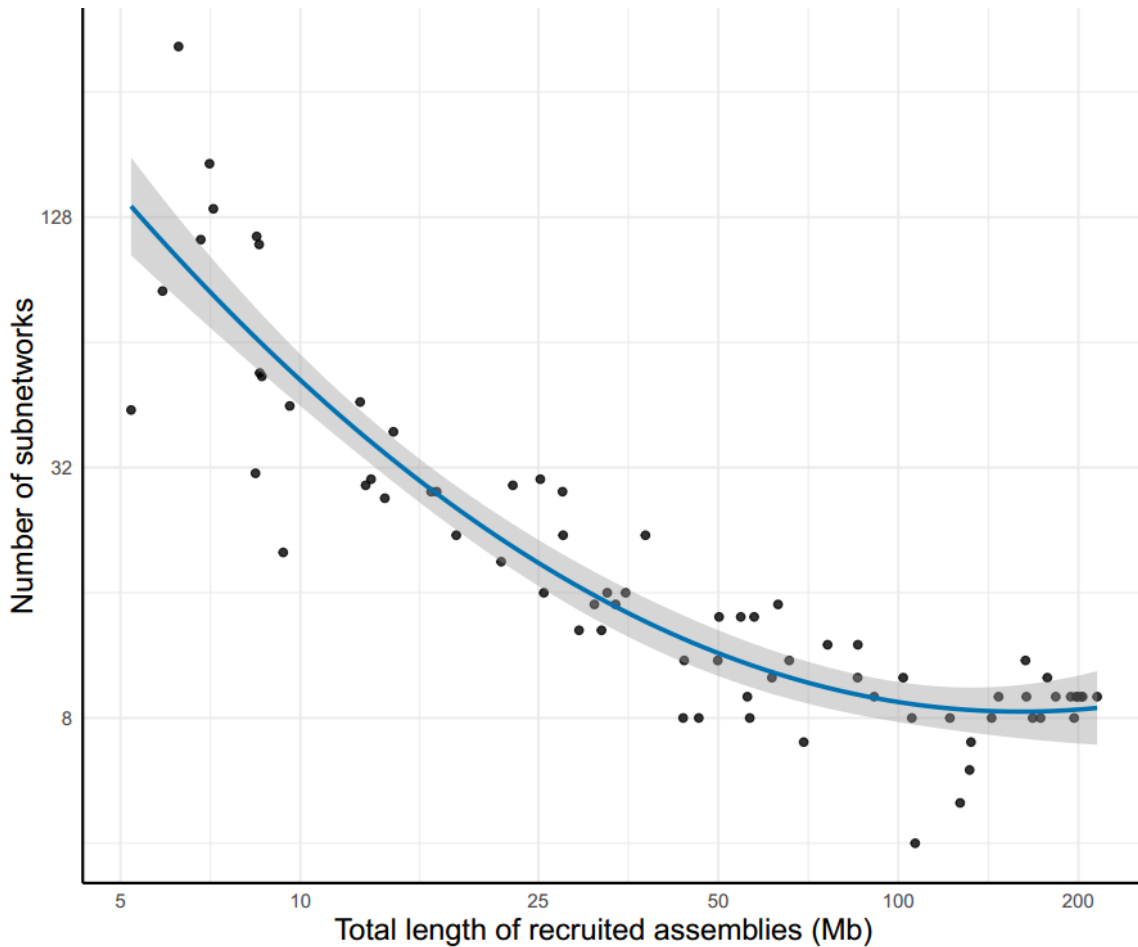


402 **Figure 2.** Subgraphs of highly variable genes in the pangenome network of 5 pathogenic *E. coli*
403 strains (manually arranged). (a) a cluster of flagellar genes. (b) a cluster containing outer
404 membrane protein-coding genes. (c) a cluster of genes responsible for biosynthesis of the O antigen.

405 (d) another cluster of O antigen-related genes. Green, blue, red nodes and edges denote reference-
406 specific, shared, and query- specific genes and gene adjacencies, respectively. Size of nodes and
407 thickness of edges indicates their weight (occurrence frequency). Numbers alongside shared genes
408 are their indexes in the representative gene set.
409



410
411 **Figure 3.** Two subgraphs of the pangenome network of *E. coli* constructed from 760 metagenomes
412 (manually arranged). (a) a cluster of flagellar genes. (b) a cluster of genes containing MGEs. Green,
413 blue, red nodes and edges denote reference-specific, shared, and query- specific genes and gene
414 adjacencies. Triangles represent MGEs. Size of nodes and thickness of edges indicates their weight
415 (occurrence frequency). Numbers alongside shared genes are their indexes in the representative
416 gene set.
417



418 **Figure 4.** Number of subnetworks in pangenome networks derived from varying sizes of recruited
419 assemblies. The x-axis indicates total length of recruited assemblies for each sub-dataset and the
420 y-axis represents the number of subnetworks in the pangenome network derived from each sub-
421 dataset. The curve was fitted for the scatters using the ‘loess’ smoothing method in R[47]. The
422 shaded area displays the 95% confidential intervals of the curve. Axes are log2-transformed.

423

424 Additional information

425 **Supplementary Figure S1.** Another cluster of genes containing MGEs, flanked by different shared
426 genes on different *E. coli* genomes (manually arranged). Green, blue, red nodes and edges denote
427 reference-specific, shared, and query-specific genes and gene adjacencies, respectively. Triangles
428 represent MGEs. Size of nodes and thickness of edges indicates their weight (occurrence
429 frequency). Numbers alongside shared genes are their indices in the representative gene set, and
430 numbers in parentheses indicate loci of these genes in the reference genome.

431 **Supplementary Figure S2.** Examples of arrangement determined by the algorithm. (a)
432 arrangements for shared nodes (blue) and reference-specific nodes (green). (b-e) arrangements for
433 query-specific nodes (red).

434 **Supplementary Table S1.** Metadata of isolate genomes used in this study.

435 **Supplementary Table S2.** Statistics for the 5 mock metagenomic datasets.

436 **Supplementary Table S3.** Tables of nodes and edges in the 5-*E. coli*-genome pangenome network
437 and the 760-metagenome pangenome network.

438 **Supplementary File S1:** Texts for, 1) steps for constructing pangenome networks, 2) steps for
439 installing the plug-in and visualizing pangenome networks in Cytoscape.

440 **Supplementary File S2:** Texts for, 1) steps for selecting query-specific nodes for arrangement, 2)
441 Comparison of the reference pangenome network (RPGN) and the query pangenome network
442 (RPGN), and 3) detailed definitions of conformity and divergence for nodes and edges.

443 **Supplementary File S3:** “5-*E. coli*-genome pangenome network.pdf”, PDF file for *E. coli*
444 pangenome network derived from five pathogenic *E.coli* strains.

445 **Supplementary File S4:** “760-metagenome pangenome network.pdf”, PDF file for *E. coli*
446 pangenome network derived from 760 genuine metagenomes.

447

448 Abbreviations

449 *E. coli*: *Escherichia coli*; LPS: lipopolysaccharide; MGEs: mobile genetic elements; *P. copri*:
450 *Prevotella copri*.

451

452 Ethics approval

453 This study has been approved by the Institutional Review Board on Bioethics and Biosafety
454 (reference number: BGI-IRB 16017).

455

456 Consent for publication

457 Not applicable.

458

459 Competing interests

460 The authors declare no competing interests.

461

462 Authors' contributions

463 J.L. conceived and directed the project. S.T. developed the plug-in. S.T., X.C., Z.Z. and Y.P.
464 developed other codes. Y.P., H.Z., J.L., D.W., S.T. and H.J. performed research. S.T. and Y.P.
465 prepared display items. J.L., H.Z., Y.P., D.W., K.K. and S.T. participated in discussion of the
466 project. Y.P., D.W., H.Z. and S.T. wrote the manuscript. All authors contributed to the revision
467 of the manuscript.

468

469 Acknowledgements

470 This study was supported by the National Natural Science Foundation of China (No.31601073).
471 We would like to express our appreciation to Dr. Liqiang Li, Dr. Ziqing Deng, Mike Huang-
472 Jingan from BGI-Shenzhen and Prof. Yue Zhang from Sichuan University, for their criticisms
473 and constructive suggestions to this study. We would like to thank Wenchen Song from BGI-
474 Shenzhen for testing the codes. We would also like to extend our gratitude to Chen Ye and Ling
475 Li from BGI-Shenzhen, who made related data publicly available, and Binge Wang and Yanmin
476 Zhao from BGI-Shenzhen for their administrative support.

477

478 References

- 479 1. Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome
480 analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the
481 microbial “pan-genome.” *Proc. Natl. Acad. Sci.* [Internet]. 2005;102:13950–5. Available from:
482 <http://www.pnas.org/cgi/doi/10.1073/pnas.0506758102>
- 483 2. Contreras-Moreira B, Vinuesa P. GET_HOMOLOGUES, a versatile software package for
484 scalable and robust microbial pangenome analysis. *Appl. Environ. Microbiol.* 2013;79:7696–
485 701.
- 486 3. Zhao Y, Wu J, Yang J, Sun S, Xiao J, Yu J. PGAP: Pan-genomes analysis pipeline.
487 *Bioinformatics.* 2012;28:416–8.
- 488 4. Cain AA, Kosara R, Gibas CJ. GenoSets: Visual Analytic Methods for Comparative
489 Genomics. *PLoS One.* 2012;7.
- 490 5. Brittnacher MJ, Fong C, Hayden HS, Jacobs MA, Radey M, Rohmer L. PGAT: A multistrain
491 analysis resource for microbial genomes. *Bioinformatics.* 2011;27:2429–30.
- 492 6. Fremez R, Faraut T, Fichant G, Gouzy J, Quentin Y. Phylogenetic exploration of bacterial
493 genomic rearrangements. *Bioinformatics.* 2007;23:1172–4.
- 494 7. Blom J, Kreis J, Spänig S, Juhre T, Bertelli C, Ernst C, et al. EDGAR 2.0: an enhanced
495 software platform for comparative gene content analyses. *Nucleic Acids Res.* 2016;44:W22–8.
- 496 8. Herbig A, Jäger G, Battke F, Nieselt K. GenomeRing: Alignment visualization based on
497 SuperGenome coordinates. *Bioinformatics.* 2012;28:7–15.
- 498 9. Pedersen TL, Nookaew I, Wayne Ussery D, Månsson M. PanViz: interactive visualization of
499 the structure of functionally annotated pangenomes. *Bioinformatics* [Internet]. 2017;33:btw761.
500 Available from: [https://academic.oup.com/bioinformatics/article-](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw761)
501 [lookup/doi/10.1093/bioinformatics/btw761](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw761)
- 502 10. Marcus S, Lee H, Schatz M, Schatz MC. SplitMEM : Graphical pan-genome analysis with
503 suffix skips *BIOINFORMATICS* SplitMEM : Graphical pan-genome analysis with suffix skips.
504 *bioArXiv.* 2014;0–7.
- 505 11. Baier U, Beller T, Ohlebusch E. Graphical pan-genome analysis with compressed suffix trees
506 and the Burrows-Wheeler transform. *Bioinformatics.* 2015;32:497–504.
- 507 12. Scholz M, Ward D V, Pasolli E, Tolio T, Zolfo M, Asnicar F, et al. Strain-level microbial
508 epidemiology and population genomics from shotgun metagenomics. *Nat. Methods* [Internet].
509 Nature Publishing Group; 2016; Available from:
510 <http://www.nature.com/doi/10.1038/nmeth.3802>

- 511 13. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics
512 pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography.
513 *Genome Res.* 2016;26:1612–25.
- 514 14. Delmont TO, Eren AM. Linking pangenomes and metagenomes: the *Prochlorococcus*
515 metapangenome. *PeerJ* [Internet]. 2018;6:e4320. Available from: <https://peerj.com/articles/4320>
- 516 15. Kim Y, Koh I, Young Lim M, Chung WH, Rho M. Pan-genome analysis of *Bacillus* for
517 microbiome profiling. *Sci. Rep.* 2017;7:1–9.
- 518 16. Farag IF, Youssef NH, Elshahed MS. Global distribution patterns and pangenomic diversity
519 of the candidate phylum “Latescibacteria” (WS3). *Appl. Environ. Microbiol.* 2017;83:1–21.
- 520 17. Cytoscape: An Open Source Platform for Complex Network Analysis and Visualization
521 [Internet]. [cited 2017 Nov 8]. Available from: <http://www.cytoscape.org/>
- 522 18. Meredith TC, Mamat U, Kaczynski Z, Lindner B, Holst O, Woodard RW. Modification of
523 lipopolysaccharide with colanic acid (M-antigen) repeats in *Escherichia coli*. *J. Biol. Chem.*
524 2007;282:7790–8.
- 525 19. Guy L, Jernberg C, Arvén Norling J, Ivarsson S, Hedenström I, Melefors Ö, et al. Adaptive
526 Mutations and Replacements of Virulence Traits in the *Escherichia coli* O104:H4 Outbreak
527 Population. *PLoS One.* 2013;8.
- 528 20. Rasko DA, Webster DR, Sahl JW, Bashir A, Boisen N, Scheutz F, et al. Origins of the *E. coli*
529 Strain Causing an Outbreak of Hemolytic–Uremic Syndrome in Germany. *N. Engl. J. Med.*
530 [Internet]. 2011;365:709–17. Available from:
531 <http://www.nejm.org/doi/abs/10.1056/NEJMoa1106920>
- 532 21. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut microbial
533 gene catalogue established by metagenomic sequencing. *Nature.* Macmillan Publishers Limited.
534 All rights reserved; 2010;464:59–65.
- 535 22. Le Chatelier E, Nielsen T, Qin J, Prifti E, Hildebrand F, Falony G, et al. Richness of human
536 gut microbiome correlates with metabolic markers. *Nature.* 2013;500:541–6.
- 537 23. Nielsen HB. Identification and assembly of genomes and genetic elements in complex
538 metagenomic samples without using reference genomes. *nbt.* 2014;2014:41–5.
- 539 24. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of reference
540 genes in the human gut microbiome. *Nat Biotech* [Internet]. 2014;advance on:834–41. Available
541 from:
542 <http://dx.doi.org/10.1038/nbt.2942>
543 [http://www.nature.com/nbt/journal/vaop/ncurrent/abs/nbt.2942.html#supplementary-](http://www.nature.com/nbt/journal/vaop/ncurrent/abs/nbt.2942.html#supplementary-information)
544 [information](http://www.nature.com/nbt/journal/v32/n8/full/nbt.2942.html?WT.ec_id=NBT)

- 545 -201408%5Cn<http://www.ncbi.nlm.nih.gov>
- 546 25. Darmon E, Leach DRF. Bacterial Genome Instability. *Microbiol. Mol. Biol. Rev.* [Internet].
547 2014;78:1–39. Available from: <http://mmbbr.asm.org/cgi/doi/10.1128/MMBR.00035-13>
- 548 26. Whitfield C, Valvano M a. Species-Wide Variation in the Escherichia coli Flagellin. *Adv.*
549 *Microb. Physiol.* 2003;35:135–246.
- 550 27. Reid SD, Selander RK, Whittam TS. Sequence diversity of flagellin (fliC) alleles in
551 pathogenic Escherichia coli. *J. Bacteriol.* 1999;181:153–60.
- 552 28. Beutin L, Delannoy S, Fach P. Sequence variations in the flagellar antigen genes
553 fliC<inf>H25</inf> and fliC<inf>H28</inf> of Escherichia coli and their use in identification
554 and characterization of enterohemorrhagic E. Coli (EHEC) O145:H25 and O145:H28. *PLoS*
555 *One.* 2015;10.
- 556 29. Heinrichs DE, Yethon JA, Whitfield C. Molecular basis for structural diversity in the core
557 regions of the lipopolysaccharides of Escherichia coli and Salmonella enterica. *Mol. Microbiol.*
558 1998. p. 221–32.
- 559 30. Iguchi A, Iyoda S, Kikuchi T, Ogura Y, Katsura K, Ohnishi M, et al. A complete view of the
560 genetic diversity of the Escherichia coli O-antigen biosynthesis gene cluster. *DNA Res.*
561 2015;22:101–7.
- 562 31. Huynen M, Snel B, Lathe W, Bork P. Predicting protein function by genomic context:
563 Quantitative evaluation and qualitative inferences. *Genome Res.* 2000;10:1204–10.
- 564 32. Delihias N. Impact of small repeat sequences on bacterial genome evolution. *Genome Biol.*
565 *Evol.* 2011;3:959–73.
- 566 33. Wang D, Li S, Guo F, Ning K, Wang L. Core-genome scaffold comparison reveals the
567 prevalence that inversion events are associated with pairs of inverted repeats. *BMC Genomics*
568 [Internet]. *BMC Genomics*; 2017;18:268. Available from:
569 <http://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-017-3655-0>
- 570 34. Doron S, Melamed S, Ofir G, Leavitt A, Lopatina A, Keren M, et al. Systematic discovery of
571 antiphage defense systems in the microbial pangenome. *Science (80-.).* 2018;1–17.
- 572 35. Serruto D, Serino L, Massignani V, Pizza M. Genome-based approaches to develop vaccines
573 against bacterial pathogens. *Vaccine.* 2009. p. 3245–50.
- 574 36. Maione D, Margarit I, Rinaudo CD, Massignani V, Scarselli M, Tettelin H, et al.
575 Identification of a Universal Group B Streptococcus Vaccine by Multiple Genome Screen.
576 2006;309:148–50.
- 577 37. Franco AA, Cheng RK, Chung GT, Wu S, Oh HB, Sears CL. Molecular evolution of the

- 578 pathogenicity island of enterotoxigenic *Bacteroides fragilis* strains. *J. Bacteriol.* 1999;
- 579 38. Sears CL, Geis AL, Housseau F. *Bacteroides fragilis* subverts mucosal biology: From
580 symbiont to colon carcinogenesis. *J. Clin. Invest.* 2014.
- 581 39. Scher JU, Szczesnak A, Longman RS, Segata N, Ubeda C, Bielski C, et al. Expansion of
582 intestinal *Prevotella copri* correlates with enhanced susceptibility to arthritis. *Elife.* 2013;
- 583 40. Laursen BS, Sørensen HP, Mortensen KK, Sperling-Petersen HU. Initiation of protein
584 synthesis in bacteria. *Microbiol. Mol. Biol. Rev.* [Internet]. 2005;69:101–23. Available from:
585 <http://www.scopus.com/inward/record.url?eid=2-s2.0-14844340954&partnerID=tZOtx3y1>
- 586 41. De Boer HA, Hui AS. Sequences within ribosome binding site affecting messenger RNA
587 translatability and method to direct ribosomes to single messenger RNA species. *Methods*
588 *Enzymol.* 1990;185:103–14.
- 589 42. Berwal SK, Sreejith RK, Pal JK. Distance between RBS and AUG plays an important role in
590 overexpression of recombinant proteins. *Anal. Biochem.* 2010;405:275–7.
- 591 43. Zhu W, Lomsadze A, Borodovsky M. Ab initio gene identification in metagenomic
592 sequences. *Nucleic Acids Res.* 2010;38.
- 593 44. Li W, Godzik A. Cd-hit: A fast program for clustering and comparing large sets of protein or
594 nucleotide sequences. *Bioinformatics.* 2006;22:1658–9.
- 595 45. Kent WJ. BLAT - The BLAST-like alignment tool. *Genome Res.* 2002;12:656–64.
- 596 46. Mende DR, Waller AS, Sunagawa S, Järvelin AI, Chan MM, Arumugam M, et al.
597 Assessment of metagenomic assembly using simulated next generation sequencing data. *PLoS*
598 *One.* 2012;7.
- 599 47. R: The R Project for Statistical Computing [Internet]. [cited 2018 Mar 6]. Available from:
600 <https://www.r-project.org/>