

- 1 Main manuscript contains the following:
2 1. Main txt (abstract, introduction, results and discussion), Line 8 – Line 579;
3 2. Online Methods, Line 588 - line 896;
4 3. Figures and Tables, Line 924 – Line 1057;
5 4. Supplementary figures, tables and notes (in Supplementary_figures_tables_notes.pdf)
6 5. Supplementary files (in Supplementary_files.zip)
7

8 **GSearch: Ultra-Fast and Scalable Microbial Genome Search**
9 **by combining Kmer Hashing with Hierarchical Navigable**
10 **Small World Graphs**

11
12 Jianshu Zhao^{1,2,7}, Jean Pierre-both^{3,7}, Luis M. Rodriguez-R^{4,5,6}, Konstantinos T.
13 Konstantinidis^{1,2,4,*}

14 ¹Center for Bioinformatics and Computational Biology, Georgia Institute of Technology,
15 Atlanta, Georgia, USA

16 ²School of Biological Sciences, Georgia Institute of Technology, Atlanta, Georgia, USA

17 ³ Université Paris-Saclay, CEA, List, Palaiseau, France

18 ⁴School of Civil and Environmental Engineering, Georgia Institute of Technology, Atlanta,
19 Georgia, USA

20 ⁵Department of Microbiology, University of Innsbruck, Innsbruck, Austria

21 ⁶Digital Science Center (DiSC), University of Innsbruck, Innsbruck, Austria

22 ⁷Those authors contribute equally

23 *Corresponding author, Konstantinos T. Konstantinidis

24 (kostas.konstantinidis@gatech.edu)

25

26

27 **Abstract**

28 Genome search and/or classification is a key step in microbiome studies and has become
29 more challenging due to the increasing number of available (reference) genomes in
30 recent years and the fact that traditional methods do not scale well with larger databases.
31 By combining a kmer hashing-based genomic distance metric (Probminhash) with a
32 graph based nearest neighbor search (NNS) algorithm (called Hierarchical Navigable
33 Small World Graphs), we developed a new program, GSearch, that is at least ten times
34 faster than alternative tools for the same purposes while maintaining high accuracy.
35 GSearch can identify/classify eight thousand query genomes against all available
36 microbial and viral genomic species within several minutes on a personal laptop, using
37 only ~6GB of memory. Further, GSearch can scale well with millions of database
38 genomes based on a database splitting strategy. Therefore, GSearch solves a major
39 bottleneck in current and future microbiome studies that require genome search and/or
40 classification.

41

42 **Keywords:** genome search, microbial genomes, MAGs, MinHash, nearest neighbor
43 search, classification, hierarchical small world graphs, HNSW

44

45

46

47

48

49 **Introduction**

50 Classifying microbial species based on either universal marker genes (e.g., 16S or
51 18S rRNA genes) or entire genomes represents a re-occurring task in environmental and
52 clinical microbiome studies. However, this task is challenging because i) whether or not
53 microbes (bacteria, fungi) and viruses form discrete population clusters (or species),
54 remains an open question ¹, and ii) the microbial species in nature are still severely under-
55 sampled by the available genomes. For instance, there are 10^{12} bacterial and fungal
56 species in nature according to a recent estimation ² and even more viral species
57 (e.g. $>10^{14}$ species). Yet, only ~17,000 bacterial species have been described and even
58 fewer (around 15,000) are represented by complete or draft genome ³. Due to the recent
59 improvements in DNA sequencing and single-cell technologies, metagenomic surveys
60 can now recover hundreds, if not thousands, of these yet-to-be-described species from
61 environmental or clinical samples ^{4,5}, filling in the gap in the described diversity mentioned
62 above. This has created a new challenge, however; that is, identifying these new
63 genomes against the exponentially increasing number of available (described) genomes
64 has become computationally intractable. Nonetheless, the recent high-throughput
65 sequencing of isolate genomes as well as metagenomic studies of natural populations
66 have shown that species may exist and be commonly circumscribed based on a 95%
67 genome-aggregate average nucleotide identity (ANI) threshold, at least for prokaryotes
68 and viruses ^{6,7}. This threshold represents convenient means in searching and identifying

69 new genomes against the already described species and determining whether or not they
70 represent novel species ⁸.

71

72 The number of curated draft or complete prokaryotic genomes has reached
73 317,542 in the newest release of the GTDB database, and 2,332,702 in the latest IMG/VR
74 database for viruses, representing 65,703 prokaryotic and 935,122 viral distinct species
75 at the 95% ANI level ^{9, 10}. Searching of query genomes against these large databases to
76 find closely-related database/reference genomes for taxonomy classification based on
77 the traditional brute-force methods, meaning, performing all vs. all searches, has become
78 impractical, even for fast searching algorithms and/or small-to-medium computer clusters.
79 For this task, faster search strategies are necessary. In addition to the searching strategy,
80 the actual algorithm used to determine overall genetic relatedness between the query and
81 the databased genomes is critical. While the traditional blast-based ANI among closely
82 related genomes at the species level, and the genome-aggregate average amino acid
83 identity (AAI) for genomes related at the genus level or above, have been proven to be
84 highly accurate for genetic relatedness estimation across microbial and viral genomes ¹¹⁻
85 ¹³, they are too slow to use when dealing with more than a few dozen genomes. Faster
86 implementations based on k-mer counting have been recently described to alleviate this
87 bottleneck such as FastANI and MASH ^{14, 15}, but these methods still do not scale with an
88 increasing number of database (or query) genomes, especially based on an all vs. all
89 search strategy. Further, defining genetic distance (or relatedness) based on kmer
90 profiles can be problematic for incomplete genomes, which are commonly recovered from
91 metagenomic surveys, and/or genomes with lots of repeats. Kmer-weighted approaches

92 are advantageous in the latter cases because repeated genomic fragments can be
93 considered when hashing but they have not been widely adopted yet ^{16, 17}. Recently, a
94 phylogeny-based approach using a handful of universal genes (n= ~100) was developed
95 to accelerate genome classification ¹⁸. However, phylogenetic replacement based on
96 concatenated universal gene tree can be memory demanding (300+ GB) and slow,
97 especially for a large number of or a few deep-branching (novel) query genomes, and this
98 approach cannot be applied to viral genomes, which lack universal genes. Further,
99 universal genes due to their essentiality, are typically under stronger (purifying) selection
100 and thus, evolve slower than the genome average. This property makes universal genes
101 appropriate for comparisons among distantly related genomes, e.g., to classify genomes
102 belong to new class or new phylum, but not the species and genus levels ^{18, 19}.

103 One of the most generally used approaches for finding closely related information
104 to a query, while circumventing an all vs. all search, is the K-Nearest Neighbor Search
105 (K-NNS). The K-NNS approach has been used for 16S rRNA gene-based classification
106 followed by a vote strategy ^{20, 21} and, more recently, for whole genome and metagenome
107 comparisons based on shared kmers ¹⁴. Approximate nearest neighbor search (ANN)
108 algorithms, such as locality-sensitive hashing (LSH) ^{22, 23}, k-dimension tree ²⁴, random
109 projection trees ²⁵, k-graph ²⁶ and proximity graph ^{27, 28} have been recently used to
110 accelerate ANN search process. Proximity graph, as implemented for example in the
111 hierarchical navigable small world graph (HNSW) ²⁹, has been shown to be one of the
112 fastest ANN search algorithms ³⁰. HNSW incrementally builds a multi-layer structure
113 consisting of a hierarchical set of proximity graphs (layers) for nested subsets of the
114 stored elements. Then, through smart neighbor selection heuristics, inserting and

115 searching the query elements in the proximity graphs can be very fast while preserving
116 high accuracy, even for highly clustered data^{27, 29}. Therefore, finding the closest genomes
117 in a database can be substantially accelerated by using HNSW.

118 Here, we describe GSearch (for Genome Search), a tool that combines the most
119 efficient nearest neighbor search approaches (HNSW) with a universal approach to
120 measure genetic relatedness among any microbial genome, including viral genomes,
121 Probminhash³¹, implemented in the Rust language for higher speed. Probminhash is
122 based on shared kmers, weighted by their abundance and normalized by total kmer size,
123 which can account for genome incompleteness of prokaryotic genomes and repeats
124 commonly found in eukaryotic and sometimes in prokaryotic genomes. Essentially,
125 Probminhash computes the normalized weighted Jaccard distance between each pair of
126 genomes and subsequently, the normalized (by total kmer size) weighted Jaccard
127 distance is used as input to build HNSW to create the graph of the database genomes.
128 Accordingly, the search of the query genome(s) against the graph to find the nearest
129 neighbors for classification purposes becomes an ultra-fast step using GSearch and can
130 be universally applied to all microbial genomes. The novelty of GSearch also includes a
131 hierarchical pipeline that involves both nucleotide-level (when query genomes have close
132 relatives at the species level) and amino-acid-level searching (when query genomes are
133 somehow novel), which provides robust classification for query genomes regardless of
134 their degree of novelty, as well as a database-splitting strategy that allows GSearch to
135 scale up well to millions of database genome sequences.

136

137 **Results**

138 *Probminhash as a robust metric for genome relatedness of prokaryotes*

139 Correlations between Probminhash distance (or we called it ProbMASH after
140 transformation) and ANI (determined by FastANI) or MASH distance showed that
141 Probminhash is robust and slightly better than MASH for determining distances among
142 bacterial genomes related at ~80% ANI, or higher, i.e., closely related genomes of the
143 same or closely-related species (Spearman $\rho=0.9643$ and 0.9640 respectively,
144 $P<0.001$, Figure S1 (a) and (b), note that for finding best hits compare to ANI, Spearman
145 rank correlation is more relevant than Pearson correlation). For moderately related
146 genomes, for which nucleotide-level ANI or distances are known to lose accuracy,
147 Probminhash was still robust compared to MASH for bacterial genomes (amino acid
148 distance or AAI), especially among genomes showing between ~52% and 95% AAI
149 (Spearman $\rho=0.90$, $P<0.01$, Supplementary Figure S2 (a) and (b)). Below ~50% AAI,
150 both Probminhash and MASH distance lose accuracy compared to AAI. However, AAI of
151 just universal genes provides a robust measurement of genetic relatedness at this level
152 of distantly related genomes¹⁹, and we show here that Probminhash distance for the
153 same set of universal genes is also robust (Spearman $\rho=0.9390$, $P<0.001$,
154 Supplementary Figure S3). Thus, for distantly related (i.e., deep-branching) query
155 genomes, e.g., their closest genome in the database is related at the order level or higher,
156 restricting the search to the universal genes can provide robust classifications.

157

158 *Graph building and search against reference prokaryotic genomes is faster than*
159 *alternative methods*

160 To build the database graph for the entire GTDB v207 database (65,703 unique, non-
161 redundant prokaryotic species) at the nucleotide level, the tohnsw module of GSearch
162 took 2.3 h on a 24-thread computing node and scaled well with increasing number of
163 threads (Figure 2 (a)). Maximum memory required for the building step was 28.3 GB. The
164 total size of written database files on disk was ~3.0 GB. There are 3 layers for the dumped
165 graph, 65180, 519 and 4 genomes for layer 0,1 and 2 respectively. The searching of query
166 genomes against this database graph, requesting best 50 neighbors for 1000 query
167 genomes, which represented different previously known as well as novel species of eight
168 bacterial phyla (see Methods for details on query genome selection), took 2.3 min
169 (database loading 6 seconds) on a 24 thread machine and also scaled well with
170 increasing number of threads (Figure 3 (a)). The memory requirement for the request
171 (search) step was only 3.0 GB for storing the entire database file in memory. To evaluate
172 the accuracy of these results, we compared the best neighbors found by GSearch with
173 brute-force FastANI and GTDB-tk. All best neighbors found by brute-force FastANI and
174 GTDB-tk for query genomes with close relatives in the database (e.g., showing >80%
175 ANI) were found by GSearch (Supplemental File 1). Top 5 neighbors were 99.4%
176 overlapping and top 10 were 96.3% overlapping between GSearch and the other two
177 methods for the testing query genomes. We also compared the speed with MASH for the
178 same kmer and sketch size and MASH dist step took 7.51 min for comparing 1000
179 genomes with database using 24 threads. The speedup compared to MASH was even
180 larger for ~8,000 query genomes. Specifically, it took 12.5 min for GSearch to find the top
181 50 best hits (Supplementary Figure S4 (a)) while MASH took 80.8 minutes on the same
182 24 thread machine. However, for a given number database genomes, speedup is

183 saturated to $\log(N)$ as the number of query genome increases, where N is the number of
184 database genomes. GSearch will be orders of magnitude faster than MASH for larger
185 species database with millions of genomic species (see also phage section). GSearch
186 query time for a given number of genomes is related to the number of database genomes
187 in a $O(\log(N))$ manner while brute-force methods are $O(N)$, and our empirical analysis is
188 consistent with the theoretical $\log(N)$ prediction (Supplementary Figure S4 (b)).

189

190 For building the amino-acid level graph for moderately related query genomes, all GTDB
191 v207 genomes are used for gene calling by FragGeneScanRs and subsequently, the
192 predicted amino acid sequences for each genome are used for the tohns module. The
193 graph building step took 1.4 h (Figure 2 (b)) with maximum memory required for the
194 building step to be 37.7 GB. The total size of written database files on disk by GSearch
195 was 5.9 GB. There were 65158, 543 and 2 genomes for layer 0,1 and 2 respectively.
196 Requesting 50 neighbors for 1000 genomes at this amino-acid level took 1.52 minutes
197 with memory requirement ~ 6.0 GB (database loading 9 seconds; Figure 3(b)). Top 5
198 neighbors were 98.9% recall with those of the brute-force MASH or blast-based AAI
199 approaches, with 97.1% overlap for the 10 top neighbors. In comparison, MASH dist took
200 5.96 min using 24 threads; for 8000 query genomes, MASH dist took 47.2 min while
201 GSearch took 5.6 min.

202

203 Finally, for most distantly related query genomes, the graph building for the universal
204 gene set follows the same logic with the amino-acid level graph mentioned above except
205 for using a smaller kmer size ($k=5$) due to the smaller kmer space of 120 universal genes

206 vs. the whole-genome level (e.g., a few thousand genes). It took 7.76 min to build the
207 database (Figure 2(c)) and 32 seconds to request 50 neighbors for 1000 queries on a 24
208 threads node (Figure 3(c)) with similar high recall to the amino-acid level search (with top
209 5 and top 10 recall, 98.2% and 97.1%, respectively).

210

211 We also evaluated the effect of genome completeness on search and classification
212 accuracy given that bacterial genomes recovered from environmental metagenomes are
213 frequently incomplete. GSearch was robust to genome incompleteness down to 50%
214 completeness level, e.g., 80% of top 10 best matches are found, while accuracy
215 decreased considerably below this level (Supplementary table S6).

216

217 *Graph database building and searching for phage and fungal genomes*

218 Graph building and requesting for phage genomes is not effective at the nucleotide level
219 because many phage genera are too diverse and do not have close relatives in the public
220 genomic database; that is, the database is too sparse. Accordingly, kmer-based methods
221 (e.g., MASH and probminhash) will often lead to imperfect graph structure for viral
222 genomes. Therefore, we build only an amino-acid level graph for viral genomes, using all
223 genes in the genome due to the lack of universal genes for viral genomes. Database
224 building took 23.895 h on a 24-thread node (Supplementary Figure S6 (a)). Request 1000
225 neighbors scales well with increasing number of threads and took about 4.4 min
226 (database load takes additional 1.9 min) using 24 threads (Supplementary Figure S6 (b)).
227 Top 10 neighbors for 1000 query phage genomes were still highly overlapping (98.32%
228 recall; Supplemental Table S1) with the brute-force MASH-based approach. For such

229 large database, GSearch is about 20X faster than the brute-force MASH (Supplementary
230 Tables S1). We also compare GSearch with a new database build method called
231 PhageCloud, which relies on manually curated genome labels (e.g., environmental
232 source) for graph database building in Neo4j database software and Dashing software
233 for distance/relatedness computation. Since PhageCloud provides only a website and
234 allows only one genome query at a time, we searched only one phage genome at a time
235 with GSearch and MASH against the same database (Gut Phage Database ³²). It took 37
236 seconds for finding the two best matches with PhageCloud while GSearch took 15
237 seconds (database loading 14 seconds, search 1 second) for the same search. MASH
238 on the other hand took 4 minutes to find the same 2 best matches. It should be noted,
239 however, that, because the database is already available (loaded) on PhageCloud's
240 website, 37 seconds is only for search and website responses (average value for 5 runs
241 on 5 different days) whereas GSearch took only 1.5 second for the same step.

242

243 Graph building for fungal genomes is slower compared to prokaryotic genomes, despite
244 the smaller number of available fungal genomes (n=9700) because the average fungal
245 genome size is much larger and kmer and sketch size are accordingly much larger (k=21,
246 s=48000). It took 2.3 h on a 24 thread node to build the nucleotide-level graph for these
247 fungal genomes. Searching step was also slower due to the larger kmer space.
248 Accordingly, it took 3.13 min for identifying 50 neighbors for 50 query fungal genomes
249 while MASH took 4.4 min. Nonetheless, recall was still very high (~99.4%) against MASH
250 and MUMMER-based ANI for the same datasets. For the amino-acid level graph, the time
251 for graph building was only 0.61 h, shorter than the corresponding prokaryotic graph due

252 to the lower coding density of fungal genomes relative to the prokaryotic genomes.
253 Identifying 50 neighbors for 50 query fungal genomes at the amino-acid level took 1.24
254 min (MASH took 2.59 min) with similar high top 5 recall (~99.7%) against brute-force
255 MASH (-a) and blastp-based AAI. Note that the difference in run time will be much larger
256 between MASH and GSearch as the number of fungal database genome increases,
257 which is similar to that of bacterial genomes

258

259 *Combining the three graphs/levels together and comparison with GTDB-tk for prokaryotic*
260 *genome classification*

261 A three-step pipeline was developed to allow the identification and classification of a
262 query genome, depending on its level of novelty compared to the database genomes
263 (Figure 4). Specifically, when the query genome does not find a match in the database
264 better than ANI > 78%, corresponding to probminhash distance 0.9850, the nucleotide-
265 level graph is abandoned, and the amino-acid level is used instead. If no match against
266 the latter graph is found above 52% AAI, corresponding to 0.9375 probminhash distance,
267 the amino-acid level is abandoned, and the universal gene graph is used instead (uAAI
268 based on universal gene below 80% indicates new order or higher taxonomic rank)(Figure
269 4). The overall running time for classifying 1000 prokaryotic genomes of varied levels of
270 taxonomic novelty on different computing platforms is showed in Table 1. On a 24 thread
271 Linux node, it took a total of 5.85 minutes while it took 19.49 minutes on an intel Core i7
272 laptop (2017 release) CPU personal laptop (6.02 minutes on the most recent ARM64
273 CPU laptop). Classifying 1000 genomes using GTDB-tk took 244.7 min on the same Linux

274 node with 24 threads (Figure 3 (d), memory requirement is ~328G) while MASH takes
275 53.7 min for 1000 genomes using 24 threads for the 3 steps.

276

277 In terms of the accuracy of the nearest neighbors found, query genomes (699 out of the
278 total 1000) that had a best match higher than 78% ANI against the GTDB database
279 genomes (i.e., a match at the same or closely related species), GSearch classified them
280 exactly the same with GTDB-tk and FastANI (Supplementary File 1, only 100/699 are
281 shown for simplicity). For the remaining 301 genomes that did not have same or closely
282 related species-level matches, for 266 of them (or 87.1%), GSearch also provided the
283 same classification with GTDB-tk but several inconsistencies were observed for 39/301
284 genomes (Supplementary Figure S5). Specifically, we noticed that for GTDB-tk, which
285 relies on RED values and tree topology, several genomes (n=14) were still classified at
286 the genus level even though AAI value against the best database genomes found was
287 below 60%, and some genomes (n=16) were still classified at the family level but not the
288 genus level even though their best AAI value was above 65%. Several genomes (n=9)
289 were classified at the order level but not family level even though their best AAI value was
290 above 52%. Therefore, high consistency was overall observed between GSearch and
291 GTDB-tk assignments, and the few differences noted were probably associated with
292 contaminated (low quality) MAGs or taxonomic inconsistencies, which was challenging to
293 assess further, and/or the peculiarities of each method. Since Probinhash distance
294 correlated well with AAI in the range of AAI values between 52% and 95%, the
295 classification results were always consistent with AAI-based classification, e.g., best
296 matches of 65% AAI or higher were classified at the same genus by GSearch and blast-

297 based AAI and best matches of $52\% < \text{AAI} < 65\%$ were typically classified at the same
298 family³³.

299

300 *Database split for large genomic species database*

301 For large databases (for example, >1 million bacterial genomes), the graph building and
302 requesting step could require a large amount of memory (due to the larger kmer space)
303 that is typically not available in a single computer node. We therefore provide a database
304 split solution for such large databases. The average database building time on each node
305 (for each piece of the database after the splitting step) scales linearly with increasing
306 nodes/processors and requires much less memory (1/n total memory compared to when
307 building in one node) (Supplementary Figure S7(a)). The searching time scales sub-
308 linearly with increasing number of nodes (Supplementary Figure S7(b)). The top 10 best
309 neighbor by splitting the database were exactly the same as the non-splitting strategy
310 (Supplementary file 2). Note that without multi-node support (e.g., run database build
311 sequentially), database build time is nearly the same with non-split strategy but memory
312 requirement is only 1/n, where n is the number of database pieces, despite the fact that
313 total request time will be larger (time*n in Supplementary Figure S7(b)). However, since
314 the request step is very fast, even for a decent number of pieces, overall runtime is still
315 short with the database split approach. The database split strategy is especially useful
316 when memory requirement is not satisfied on host machine for larger genomic species
317 database (e.g., millions of genomes).

318

319 **Discussion**

320 A popular way to assess genetic relatedness among genomes is ANI, which
321 corresponds well to both 16S/18S rRNA gene identity and DNA-DNA hybridization values,
322 the golden standards of fungal and prokaryotic taxonomies ¹¹. As the number of available
323 microbial genomes has grown at an unprecedented speed recently (e.g., there are 30%
324 more (new) species in GTDB v202 (2020) vs. v207 (2022), and the number of bacterial
325 species genomes alone is expected to surpass 1 million in the near future), the traditional
326 way that blast-based ANI or faster kmer-based implementations (e.g., FastANI or MASH)
327 are applied as an all vs. all search strategy (brute-force) does not scale because the
328 running time grows linearly with increasing number of query genomes and/or genomes in
329 the database. Phylogenetic approaches based on quick (approximate) maximum
330 likelihood algorithms and a handful of universal genes as implemented -for example- in
331 GTDB-Tk could be faster than brute-force approaches but are often not precise and
332 require a large amount of memory for querying step ^{18, 34} while the database building step
333 could take several weeks of run time because the underlying multiple sequence alignment
334 of the database genomes is computationally intensive. Further, approaches that rely on
335 k-medoid clustering to avoid all vs. all comparisons could be sometimes trapped into local
336 minima because of arbitrary partitioning of database genomes into clusters, a known
337 limitation of these methods ¹⁹. Our GSearch software effectively circumvents these
338 limitations by combining a new kmer hashing-based (Probminhash) algorithm for fast
339 computation of genetic relatedness among genomes with a graph based nearest neighbor
340 search algorithm. Accordingly, GSearch is at least an order of magnitude faster than
341 alternative approaches for the same purposes. Note that GSearch could also be applied
342 to whole metagenome comparisons and identification of the most similar metagenomes

343 in a series because ProbMinhash can estimate metagenomic distance in a similar way to
344 genomes.

345 To the best of our knowledge, no current tool can efficiently search very large
346 genome databases. GSearch is able to handle a million microbial genomes on a small-
347 to-average computer cluster since the dumped database file size is proportional to the
348 total number of genomes in database for fixed sketch size and graph parameters.
349 Specifically, with a million genome species, dumped file size (amino acid) will be
350 $5.9G \times 20 = 118$ GB, a modest computational requirement for current computer clusters.
351 Further, due to the nature of graph based NNS algorithm, there is no need to build the
352 entire database at once, but the database can split it into smaller pieces and thus, a
353 separate graph database be built for each piece as exemplified above and depending on
354 the computational resources available. For a modern laptop with 16 GB memory, a
355 database on one million species can be split into 10 pieces, so dumped file for each piece
356 will be only 11.8 GB, which can be loaded into memory, and then collect the results from
357 each piece within an approximate total running time of 30 minutes (assume each part will
358 be 3 minutes) (Supplementary table S7). With this logic, a computing node with 24 threads
359 and 256 GB of memory available can easily deal with 20 million bacterial genomes. This
360 represents a major improvement compared to existing tools for the same purposes.

361 It is also important to note that we could seemingly replace Probminhash with
362 another relatedness algorithm should such an algorithm become available and has
363 advantages in terms of speed and/or precision. Related to this, ANI as currently
364 implemented -for instance- in FastANI is not appropriate for this function because it is not
365 metric (that is, for the FastANI distances calculated among three genomes A, B, and C,

366 (A,B) + (B, C) is not necessary larger than (A,C), especially for genomes related at the
367 phylum level). To solve this “metric” problem, a norm adjusted proximity graph (NAPG)
368 was proposed based on inner product and it shows improvements in terms of both speed
369 and recall ³⁵. This could be another direction for further improving the speed and recall of
370 GSearch and/or the use of other metrics in place of Prominhash distances. In the
371 meanwhile, Probminhash was used in GSearch because it is metric, which ensures
372 neighbor diversity when building the graph, but also equally applicable to any microbial
373 genome, including viral genomes, in addition to its advantages for kmer weighting and
374 normalization mentioned above.

375 Another distinguishing aspect of GSearch (tohns module) is the speed and
376 flexibility in building reference databases. Indeed, users could build reference databases
377 (graphs) for any number and type (e.g., prokaryotic vs. viral) of genomes, up to several
378 millions of genomes. The high efficiency in building graphs allows users to also test and
379 optimize the key parameters of the graph, the M and ef_construct parameters. For any
380 given database size, M and ef_construct determine the quality of the graph and graph
381 build speed. Small M and ef_construct may lead to frequent traps in local minima and
382 thus, low recall while large M and ef_construct may lead to slow speed without
383 proportional improvement in recall (Supplemental Table S2). Therefore, there is a tradeoff
384 between accuracy and speed that should be evaluated first. However, for most users this
385 task would not be necessary because they will work with pre-built databases such as
386 those provided here. Further, the search step against these pre-build databases with
387 query genomes of known taxonomy for evaluating recall and tradeoffs can be performed,
388 within minutes, on any modern laptop with 5-6 GB of memory (Table 1).

389 Kmer-based methods for genetic relatedness estimation such as Probminhash
390 have lower accuracy between moderately-to-distantly related genomes compare to
391 alignment-based tools (see supplement Note 4 for further discussion). Our empirical
392 evaluation showed that this relatedness level, for nucleotide searches, is around 78% ANI
393 and 52% AAI for the amino-acid searches (e.g., probminhash distances do not correlate
394 well with blast-based ANI and AAI at these levels). To circumvent this limitation, we
395 designed a 3-step framework as part of GSearch to classify bacterial genomes that show
396 different levels of novelty compared to the database genomes, with high accuracy. This
397 framework included a search at the universal gene level for deep-branching genomes
398 that are novel at the phylum level (e.g., showing <52% AAI), for which searching at the
399 entire proteome level is less accurate. Recently, methods that employ kmers that allow
400 mismatches, that is, spaced kmers ³⁶, have shown promise in accurately estimating
401 genomic relatedness even among distantly related genomes with gains in speed and has
402 already been applied in classification. To apply spaced kmer to entire genomes, the
403 recently developed “tensor sketch” approaches could be explored in the future to simplify
404 the pipeline for bacterial and viral genomes ³⁷. In the meanwhile, the probminhash
405 approach, essentially a Jaccard distance estimation via MinHash-based analysis of kmers, is
406 highly efficiently, and, importantly, can effectively deal with incomplete genomes or
407 genomes of (drastically) different length, an known limitation of MASH-based methods ³⁸.
408 Comparing genomes of different length is not uncommon, e.g., bacterial genome size can
409 differ by more than two-fold, as can be the case between MAGs of different level of
410 completeness or when searching a short sequence (e.g., a bacteriophage genome)
411 against a large genome collection (e.g., whole viral genome database). Probminhash is

412 more robust with genomes of unknown completeness by weighting completed genomes
413 more due to the kmer normalization step, and our own analysis showed that it is robust
414 down to 50% completeness level (Supplemental Table S6), which is also the most
415 commonly used standard for selecting MAGs of sufficient/high quality ³⁹.

416 In general, the genome relatedness estimated, or best database matching genome
417 identified by GSearch were highly consistent with Blast-based AAI results or phylogenetic
418 placement of the genome using GTDB-tk, particularly for query genomes with close
419 relatives in the database related at the species or genus level (Supplementary File 1,
420 Supplementary Figure S5). For more distantly related query genomes relative to database
421 genomes, classification results of GSearch showed some differences with GTDB-tk.
422 These differences were not always possible to assess further for the most correct genome
423 placement but could be due, at least partly, to the incompleteness and/or contamination
424 of query or/and database genomes, which renders the resulting concatenated alignment
425 of a few universal genes used by GTDB-tk unreliable ⁴⁰ (and it is a few amino-acid
426 positions per gene that are used in the final alignment). In contrast, the AAI and
427 Probminhash approaches should be more robust to changes of a small number of genes
428 because the entire proteome is considered ¹⁵.

429 Graph-based NNS methods achieve good performance compared to tree based
430 and locality-sensitive hashing (LSH) methods. Building a HNSW graph relies on proximity
431 of database element; so, if the distances among database elements, in our case
432 genomes, cannot be effectively estimated via hashing algorithms, the navigation in graph
433 will be less efficient (e.g., get trapped in local minima) because the edges to choose from
434 will not be accurate estimations of the targeted genomes they represent. This is especially

435 true for highly sparse/distantly related and diverse dataset, like the viral genome
436 database, e.g., two phage genomes could often share very little genomic information
437 (kmers) in the current dataset. This is confirmed by our own results when using
438 nucleotide-level search to build the viral graph. Hence, the amino acid level will be much
439 more robust for viral genomes and is the recommended level to use. Finally, the HNSW
440 graph, and graph-based K-NNS in general, can be further improved by adding shortcut
441 edges and maintaining a dynamic list of candidates, compared to a fixed list of candidates
442 by default ⁴¹. Graph reordering, a cache optimization that works by placing neighboring
443 nodes in consecutive (or near-consecutive) memory locations, can also be applied to
444 improve the speed of HNSW ⁴². Another new direction for graph based NNS will be using
445 Graphics Processing Unit (GPU) instead of CPU because GPUs are more efficient in
446 handling matrix computations and machine learning tasks ⁴³. We will explore these
447 options in future version of GSearch.

448

449 To summarize, GSearch, based on Probminhash and HNSW, solves a major
450 current challenge in classification of microbial genomes, especially given the exponential
451 increase in the number of newly sequenced genomes due to its efficiency and scalability.
452 GSearch will serve the entire microbial sciences for several years to come since it can be
453 applied to fungal, bacterial and viral genomes and will accelerate the process to find new
454 biological knowledge.

455

456 **Data availability**

457 All the mentioned pre-built database for bacteria, fungi and phage genomes can be found
458 at: <http://enve-omics.ce.gatech.edu/data/gsearch>

459

460 **Author Contribution**

461 J.Z, L.M and K.K designed the work, J.Z and J.P-B wrote the code (Genome part and
462 algorithm part respectively), J.P-B implemented the Rust libraries of Kmerutils,
463 Probinhash and Hnswlib-rs. J.Z and K.K wrote the paper. J.Z did the analysis and
464 benchmark.

465

466 **Acknowledgment**

467 This work was supported by the US National Science Foundation (Award No 1759831
468 and 2129823 to KTK). We want to thank PACE (Partnership for Advanced Computing
469 Environment) at Georgia Tech for providing computing resources. We want to thank Kenji
470 Gerhardt for helpful discussion on benchmarking against traditional ANI/AAI based tools
471 and Chirag Jain for discussion on the graph based nearest neighbor search.

472

473 **Reference**

- 474 1. Bobay, L.-M. & Ochman, H. Biological species in the viral world. *Proceedings of the*
475 *National Academy of Sciences* **115**, 6040-6045 (2018).
- 476 2. Locey, K.J. & Lennon, J.T. Scaling laws predict global microbial diversity. *Proceedings of*
477 *the National Academy of Sciences* **113**, 5970-5975 (2016).
- 478 3. Federhen, S. Type material in the NCBI Taxonomy Database. *Nucleic Acids Research*
479 **43**, D1086-D1098 (2014).
- 480 4. Almeida, A. et al. A unified catalog of 204,938 reference genomes from the human gut
481 microbiome. *Nature Biotechnology* **39**, 105-114 (2021).
- 482 5. Nayfach, S. et al. A genomic catalog of Earth's microbiomes. *Nature Biotechnology* **39**,
483 499-509 (2021).
- 484 6. Caro-Quintero, A. & Konstantinidis, K.T. Bacterial species may exist, metagenomics
485 reveal. *Environmental microbiology* **14**, 347-355 (2012).
- 486 7. Deng, L. et al. Viral tagging reveals discrete populations in *Synechococcus* viral genome
487 sequence space. *Nature* **513**, 242-245 (2014).

- 488 8. Rodriguez-R, L.M., Jain, C., Conrad, R.E., Aluru, S. & Konstantinidis, K.T. Reply to: “Re-
489 evaluating the evidence for a universal genetic boundary among microbial species”.
490 *Nature Communications* **12**, 4060 (2021).
- 491 9. Parks, D.H. et al. GTDB: an ongoing census of bacterial and archaeal diversity through a
492 phylogenetically consistent, rank normalized and complete genome-based taxonomy.
493 *Nucleic Acids Research* (2021).
- 494 10. Roux, S. et al. IMG/VR v3: an integrated ecological and evolutionary framework for
495 interrogating genomes of uncultivated viruses. *Nucleic acids research* **49**, D764-D775
496 (2021).
- 497 11. Konstantinidis, K.T. & Tiedje, J.M. Genomic insights that advance the species definition
498 for prokaryotes. *Proceedings of the National Academy of Sciences* **102**, 2567-2572
499 (2005).
- 500 12. Konstantinidis, K.T. & Tiedje, J.M. Towards a Genome-Based Taxonomy for
501 Prokaryotes. *Journal of Bacteriology* **187**, 6258-6264 (2005).
- 502 13. Goris, J. et al. DNA–DNA hybridization values and their relationship to whole-genome
503 sequence similarities. *International Journal of Systematic and Evolutionary Microbiology*
504 **57**, 81-91 (2007).
- 505 14. Ondov, B.D. et al. Mash: fast genome and metagenome distance estimation using
506 MinHash. *Genome Biology* **17**, 132 (2016).
- 507 15. Jain, C., Rodriguez-R, L.M., Phillippy, A.M., Konstantinidis, K.T. & Aluru, S. High
508 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.
509 *Nature Communications* **9**, 5114 (2018).
- 510 16. Brown, C.T. & Irber, L. sourmash: a library for MinHash sketching of DNA. *Journal of*
511 *Open Source Software* **1**, 27 (2016).
- 512 17. Bovee, R. & Greenfield, N. Finch: a tool adding dynamic abundance filtering to genomic
513 MinHashing. *Journal of Open Source Software* **3**, 505 (2018).
- 514 18. Chaumeil, P.-A., Mussig, A.J., Hugenholtz, P. & Parks, D.H. GTDB-Tk: a toolkit to
515 classify genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925-1927
516 (2019).
- 517 19. Rodriguez-R, L.M. et al. The Microbial Genomes Atlas (MiGA) webserver: taxonomic
518 and gene diversity analysis of Archaea and Bacteria at the whole genome level. *Nucleic*
519 *Acids Research* **46**, W282-W288 (2018).
- 520 20. Wang, Q., Garrity, G.M., Tiedje, J.M. & Cole, J.R. Naïve Bayesian Classifier for Rapid
521 Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Applied and*
522 *Environmental Microbiology* **73**, 5261-5267 (2007).
- 523 21. Schloss, P.D. et al. Introducing mothur: Open-Source, Platform-Independent,
524 Community-Supported Software for Describing and Comparing Microbial Communities.
525 *Applied and Environmental Microbiology* **75**, 7537-7541 (2009).
- 526 22. Indyk, P. & Motwani, R. in Proceedings of the thirtieth annual ACM symposium on
527 Theory of computing 604-613 (1998).
- 528 23. Gionis, A., Indyk, P. & Motwani, R. in Vldb, Vol. 99 518-529 (1999).
- 529 24. Bentley, J.L. Multidimensional binary search trees used for associative searching.
530 *Communications of the ACM* **18**, 509-517 (1975).
- 531 25. Dasgupta, S. & Sinha, K. in Proceedings of the 26th Annual Conference on Learning
532 Theory, Vol. 30. (eds. S.-S. Shai & S. Ingo) 317--337 (PMLR, Proceedings of Machine
533 Learning Research; 2013).
- 534 26. Dong, W., Moses, C. & Li, K. in Proceedings of the 20th international conference on
535 World wide web 577-586 (2011).
- 536 27. Malkov, Y., Ponomarenko, A., Logvinov, A. & Krylov, V. Approximate nearest neighbor
537 algorithm based on navigable small world graphs. *Information Systems* **45**, 61-68
538 (2014).

- 539 28. Fu, C., Xiang, C., Wang, C. & Cai, D. Fast approximate nearest neighbor search with the
540 navigating spreading-out graph. *arXiv preprint arXiv:1707.00143* (2017).
- 541 29. Malkov, Y.A. & Yashunin, D.A. Efficient and Robust Approximate Nearest Neighbor
542 Search Using Hierarchical Navigable Small World Graphs. *IEEE Transactions on Pattern
543 Analysis and Machine Intelligence* **42**, 824-836 (2020).
- 544 30. Aumüller, M., Bernhardsson, E. & Faithfull, A. ANN-Benchmarks: A benchmarking tool
545 for approximate nearest neighbor algorithms. *Information Systems* **87**, 101374 (2020).
- 546 31. Ertl, O. ProbMinHash – A Class of Locality-Sensitive Hash Algorithms for the
547 (Probability) Jaccard Similarity. *IEEE Transactions on Knowledge and Data Engineering*,
548 1-1 (2020).
- 549 32. Camarillo-Guerrero, L.F., Almeida, A., Rangel-Pineros, G., Finn, R.D. & Lawley, T.D.
550 Massive expansion of human gut bacteriophage diversity. *Cell* **184**, 1098-1109.e1099
551 (2021).
- 552 33. Konstantinidis, K.T., Rosselló-Móra, R. & Amann, R. Uncultivated microbes in need of
553 their own taxonomy. *The ISME Journal* **11**, 2399-2406 (2017).
- 554 34. Chaumeil, P.-A., Mussig, A.J., Hugenholtz, P. & Parks, D.H. GTDB-Tk v2: memory
555 friendly classification with the Genome Taxonomy Database. *bioRxiv*,
556 2022.2007.2011.499641 (2022).
- 557 35. Tan, S. et al. in Proceedings of the 27th ACM SIGKDD Conference on Knowledge
558 Discovery & Data Mining 1552–1560 (Association for Computing Machinery, Virtual
559 Event, Singapore; 2021).
- 560 36. Břinda, K., Sykulski, M. & Kucherov, G. Spaced seeds improve k-mer-based
561 metagenomic classification. *Bioinformatics* **31**, 3584-3592 (2015).
- 562 37. Joudaki, A., Rättsch, G. & Kahles, A. Fast Alignment-Free Similarity Estimation By
563 Tensor Sketching. *bioRxiv* (2021).
- 564 38. Koslicki, D. & Zabeti, H. Improving MinHash via the containment index with applications
565 to metagenomic analysis. *Applied Mathematics and Computation* **354**, 206-215 (2019).
- 566 39. Bowers, R.M. et al. Minimum information about a single amplified genome (MISAG) and
567 a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nature
568 Biotechnology* **35**, 725-731 (2017).
- 569 40. Tan, G. et al. Current Methods for Automated Filtering of Multiple Sequence Alignments
570 Frequently Worsen Single-Gene Phylogenetic Inference. *Systematic Biology* **64**, 778-
571 791 (2015).
- 572 41. Prokhorenkova, L. & Shekhovtsov, A. in Proceedings of the 37th International
573 Conference on Machine Learning, Vol. 119. (eds. D. Hal, III & S. Aarti) 7803--7813
574 (PMLR, Proceedings of Machine Learning Research; 2020).
- 575 42. Coleman, B., Segarra, S., Shrivastava, A. & Smola, A. Graph Reordering for Cache-
576 Efficient Near Neighbor Search. *arXiv preprint arXiv:2104.03221* (2021).
- 577 43. Groh, F., Ruppert, L., Wieschollek, P. & Lensch, H. GGNN: Graph-based GPU Nearest
578 Neighbor Search. *IEEE Transactions on Big Data*, 1-1 (2022).

579
580
581
582
583
584
585
586
587

Methods and Material/Online Methods

588 Briefly, GSearch is composed of the following steps. Initially, the genetic relatedness
589 among a collection of database genomes is determined based on the probability MinHash
590 algorithm (or Probminhash), which computes the normalized weighted Jaccard distance
591 using the probminhash3a algorithm implemented in the probminhash ¹. The normalized
592 weighted Jaccard distances are then used as input for building HNSW graphs (note that
593 a distance computation is required only when that genome pair is required for graph
594 building, thus GSearch avoids all vs. all distance computations). Genomes are
595 subsequently recursively added as the nearest neighbors of each node in the built graph
596 file with the same distance computation procedure. The built graph database file is stored
597 on disk. Query genomes are then searched against graph database and subsequently,
598 best neighbors are returned for classification. In this process, the best neighbor (or
599 neighbors) is also identified based on the smallest normalized weighted Jaccard distance
600 obtained.

601

602 *Probminhash*

603 MASH is a hashing-based algorithm based on MinHash ², which is very efficient for
604 comparing genome/metagenome overall similarity ³. MASH distances represent a kmer-
605 based overall overlap between sequences according to a minimal evolutionary model.
606 Essentially, MASH distance is the Jaccard similarity value of kmer shared between
607 sequence sets A and B. However, MASH, and similar MinHash-based tools, have several
608 limitations; most notably, the loss of k-mer frequency information (only presence/absence
609 of kmer is counted) and the impact of relative set size (e.g., completeness level of a
610 genome) on the Jaccard similarity estimates ^{3, 4} Although some recent MinHash

611 implementations address these limitations (e.g., the over-sketching and track-abundance
612 methods of the MinHash-based tools finch , sourmash' or FracMinHash and
613 HyperLogLog) ⁵⁻⁸, they do not utilize the frequencies of all observed k-mers in generating
614 the kmer-profile (sketch) for a given sequence set. More recently, in the HULK software,
615 consistent weighted sampling (D²histosketch, P-MinHash algorithm was proposed for J_p)
616 ⁹ was utilized to incorporate k-mer frequency information when estimating weighted and
617 standard Jaccard similarity, which effectively addresses these limitations mentioned
618 above ¹⁰. Notably, the hash algorithm (P-MinHash) used in D²histosketch could be further
619 optimized to achieve a time complexity below $O(nm)$ (where m denotes the signature size
620 and n is the number of elements with nonzero weight in two sequence sets), further
621 improving the performance of applications such as HULK. Motivated by the
622 SuperMinHash for conventional Jaccard similarity estimation ¹¹ and BagMinHash
623 algorithm for weighted Jaccard similarity estimation ¹², probminhash (probminhash 3(a)
624 and 4 algorithm) is orders of magnitude faster than the original algorithm P-MinHash
625 proposed in D²histosketch ¹. Probminhash estimates the Jaccard probability J_p index, and
626 $1 - J_p$ is indeed a metric on probability distributions and is Pareto optimal (Supplementary
627 Note 1) ^{1, 13}. Therefore, we reimplemented the Probminhash algorithm in Rust to estimate
628 genetic relatedness between any two genomes based on normalized (weighted) Jaccard
629 distances according to the original ProbMinHash paper ¹ (Supplementary Note 1) . The
630 Rust reimplementation of Probminhash can be found at: [https://github.com/jean-](https://github.com/jean-pierreBoth/probminhash)
631 [pierreBoth/probminhash](https://github.com/jean-pierreBoth/probminhash). Two important parameters of Probminhash are the sketch size
632 and kmer size. Similar to MinHash sketches, Probminhash sketches are also shared

633 hashes from hashed kmer set by taking into account of kmer weights (See Figure 1 of
634 MASH paper). Time complexity analysis for ProbMinHash is in Supplementary Note 3.

635 To benchmark probminhash against MASH, we use the same sketch size
636 ($s=12000$) and kmer size ($k=16$) for bacterial genomes at the nucleotide level and kmer
637 size ($k=7$) at the amino acid level for both database building and searching. For fungal
638 genomes a larger sketch size (48000) was used due to much larger genome size Details
639 of kmer choosing logic can be found in Supplementary Note 2. For graph search results,
640 we also perform the same transformation of MASH distance from normalized weighted
641 Jaccard distance to probMASH distance for convenience to compare with ANI based

$$probMASH = -\frac{1}{k} \ln \frac{2 * J_p}{J_p + 1}$$

642 methods.

643

644 *Hierarchical Navigable Small World Graphs (HNSW)*

645 Generally, the framework of graph-based ANN search algorithm (here HNSW) can be
646 summarized as the following two steps: 1) build a proximity graph (HNSW) where each
647 node represents a database vector. Each database vector will connect with a few of its
648 neighbors while maintaining small world property in each layer of HNSW. 2) Given a query
649 vector (or sequence, kmer profile in our case), perform a greedy search on the proximity
650 graph by comparing the query vector with database vectors under the searching
651 measures (e.g., cosine similarity or L2 similarity, in our case probminhash distance).
652 Then, the most similar candidates are returned as outputs. The key point for these two-
653 step methods is step 1, to construct a high-quality index graph, which provides a proper
654 balance between the searching efficiency and effectiveness. To guarantee the searching

655 efficiency, the degree (number of maximum allowed neighbors, denoted as M) of each
656 node is usually restricted to a small number (normally 20~200) while width of search for
657 neighbor during inserting (denoted as $ef_construct$) is usually a larger number (above
658 1000) to increase the chance to find best M neighbors by increasing the diversity of
659 neighbors due to the larger number of them. Building graph and searching query against
660 the graph follow very similar greedy search procedures except that there is an extra
661 reverse updating of neighbors list for each vector when inserting database vector
662 (building), one by one, into the existing graph (Figure 1 (a)). The first phase of the
663 insertion/building process starts from the top layer by greedily traversing the graph in
664 order to find maximum M closest neighbors to the inserted element P in the layer by doing
665 $ef_construct$ times search (Figure 1 (a)). After that, the algorithm continues the search
666 from the next layer using the closest neighbor found from the previous layer as entry
667 point, and the process repeats until to the bottom layer. Closest neighbors at each layer
668 are found by a greedy and heuristic search algorithm (Figure 1 (b) and (c)). For building,
669 after searches are finished at the bottom layer for each inserted element, a reverse update
670 step will be performed to update the neighbor list of each node in the existing graph while
671 for searching this is not needed. The overall database building time complexity is
672 $O(N \cdot \log(N))$, where N is the number of nodes in the graph. For searching, since there is
673 no need to reverse update best neighbor list for each node in the graph, time complexity
674 is (only) $O(\log(N))$ (See Supplementary Note 3). Theoretical guarantee of graph-based
675 algorithm can be found in Supplementary Note 5.

676

677 Program *implementation details in Rust*

678 We reimplemented the original hnsplib library written in C++ using the Rust programming
679 language for its memory safety and thread use efficiency ¹¹. To benchmark our
680 reimplementation of hnsplib, we followed standard ANN benchmark procedures using
681 two popular testing datasets (MINST and SIFT1M) based on their Euclidean distance¹⁴.
682 Our results showed that, for the MINST fashion dataset (784 dimensions, 60,000 vectors),
683 recall for top 100 neighbors of 10,000 query vectors is greater than 98% for a smaller
684 number of M and ef_construct, and even higher recall rate (99.86%) for a medium M and
685 ef_construct while query speed is not compromised (Supplemental Table S2). For the
686 SIFT1M dataset (128 dimensions, 1,000,000 vectors), recall for top 100 neighbors of
687 10,000 query vectors was 99.77% for a medium M and ef_construct (Supplemental Table
688 S3 and S4). The Rust package hnsplib-rs can be found at: [https://github.com/jean-](https://github.com/jean-pierreBoth/hnsplib-rs)
689 [pierreBoth/hnsplib-rs](https://github.com/jean-pierreBoth/hnsplib-rs). For each genomic database, we chose M and ef_construct
690 experimentally, by gradually increasing M and ef_construct while monitoring query speed
691 and recall, similar to what is shown in Supplementary Table S2 for MNIST dataset. We
692 stopped the assessment when there was only a marginal increase in accuracy but decent
693 decrease in speed. To leverage between recall and speed, we use M=128 and
694 ef_search=1600 for graph building for GTDB database fungal database while M=128,
695 ef_search=3200 for phage database. There are 2 modules in total: tohnsw and request.
696 Tohnsw is to build graph by gradually inserting genomes into graph while request is to
697 query new genomes against the graph database built in the tohnsw step. Tohnsw starts
698 from reading database genomes and generating kmer profile and sketches for distance
699 calculation. By selecting a random genome as the first genome to insert to the graph,
700 tohnsw module gradually add genomes to existing graph file following HNSW constructing

701 rules mentioned above by computing Probminhash distance between genomes.
702 Whenever a genome is going to be inserted into the existing graph, each genome in the
703 graph will hold a list that store the M closest neighbors/genomes that are linked to itself
704 and the distance to these neighbors. Then the distances of this genome with the nearest
705 neighbors (M) of entry genome in this layer will be computed/searched (ef_construct
706 times) using Probminhash3a algorithm and the smallest distance of the neighbor
707 genomes will be the new entry genome. This process will be repeated until the nearest
708 genomes ($\leq M$) in the layer are found and subsequently, the program will go to the layer
709 below it using the genome that was represented by the nearest genome in the above
710 layer as new entry genome in the new layer. The search layer algorithm is repeated until
711 to the bottom layer is reached/analyzed. In contrast to the default settings in the original
712 hnsplib, we allow the two parameters of neighbor selecting heuristics, *extendCandidates*
713 to be true and *keepPrunedConnections* to be false because our genomic data is
714 extremely clustered and there is no need to fix the number of connections per element
715 considering the maximum connection allowed. Request module will load the graph
716 database and then search query genomes against it to return best neighbors of each
717 query following exact the same procedure with building step without updating the
718 database. Both tohns and request module are paralleled for high performance (see
719 Supplementary Note 6). The GSearch software can be found here:
720 <https://github.com/jean-pierreBoth/archaea> or here:
721 https://gitlab.com/Jianshu_Zhao/archaea. GSearch relies on Kmerutils
722 (<https://github.com/jean-pierreBoth/kmerutils>), which is a Rust package to manipulate

723 genomic fasta files including kmer string compression, kmer counting, filtering using
724 cuckoo filter et.al.

725 Installation guide, manual and pre-built binaries can also be found. We provide static
726 binaries on the release page for major platforms such as Linux and MacOS, with support
727 for different CPU structures, e.g. Intel x86_64 or ARM64. GSearch program can be run
728 like this : 1) Build a graph database, which can be done running the following command:
729 `tohnsw -d ./GTDB_r207 -k 16 -s 12000 -n 128 --ef 1600`; 2) Request neighbors of query
730 genomes: `request -b . -r ../query_folder -n 50 (--aa)`.

731

732 *Prokaryotic classification pipeline*

733 The amino-acid level graph showed that closest neighbors were found, with high
734 recall, when the query shared at least 52% AAI to its best neighbor. For more divergent
735 genomes, showing lower than 52% AAI equivalent, whole-genome amino-acid level graph
736 loses accuracy and we had to switch to universal, single-copy protein-coding genes. For
737 the nucleotide-level graph, we used kmer=16 for bacteria and archaea to have high
738 specificity for closely related database genomes (95% ANI to each other in GTDB
739 database). For building the whole-genome amino-acid graph, we used k=7 to have the
740 best specificity without *compromising* sensitivity, which is also consistent with previous
741 research on amino acid sequences classification based on kmers¹⁵. For building graph
742 based on universal gene set, we use k=5 because of much smaller total amino acid size.
743 For details on the range of kmer that could be used for bacteria genome and proteome,
744 bacteriophage genome and proteome, see Supplemental Notes 2.

745 The proteome of each genome was predicted by FragGeneScanRs for
746 performance purpose compare to Prodigal despite small loss in precision (Supplementary
747 Table S5) ¹⁶. Hmsearch in the hmmer software ¹⁶ was used to extract universal gene
748 collection for bacteria and archaea genome for the universal gene graph. Note that for
749 phage genomes, this last step was not used because there is no universal single copy
750 genes for viral genomes. Evaluation of the speed and memory requiremetns for all steps
751 mentioned above were performed on a RHEL (Red Hat Enterprise Linux) v7.9 with 2.70
752 GHz Intel(R) Xeon(R) Gold 6226 CPU. Unless noted otherwise, all 24 threads of the CPU
753 are available by default.

754

755 *Distributed implementation and database splitting*

756 To accommodate the increasing number of genomes that become available at an
757 unprecedented speed in recent years and will soon reach 1 million or more, we provide
758 an option to randomly split the database into a given number of pieces and build graph
759 database separately for each piece. In the end, all best neighbors returned from each
760 piece will be pooled and sorted by distance to have a new best K neighbor collection
761 returned to the user for each query genome. We hereby prove that in terms of requesting
762 top K best neighbors, the database split strategy is equivalent to non-split database
763 strategy as long as the requested best neighbors for each database piece is larger than
764 or equals to requested best neighbors in the non-split strategy. The underlying reason is
765 that the best neighbors globally are also the best locally ¹⁷. The database split and request
766 will be done sequentially, on one node, without multi-node support. For now, we split
767 GTDB database in to 5 pieces for testing purposes. In theory, a large database can be

768 split into any pieces as long as each piece can be used to build HNSW. In practice, a
769 reasonable way is to split so that memory requirement for each piece is equal or smaller
770 than the total memory of host machine. The database split idea has been used in several
771 graph-based larger scale (e.g., billions) nearest neighbor search tasks in industry ^{17, 18}.

772

773 *Species database and testing genomes for benchmarking and recall*

774 GTDB version 207 was used to build database for bacteria and archaea genome species
775 ¹⁹. It appears that virtually all metagenome-assembled dsDNA viral populations form
776 discrete genotypic clusters/species and can be appropriately delineated using a $\geq 95\%$
777 genome-wide ANI cut-off ²⁰. The IMGVR database version 3, with species representatives
778 at a $\geq 95\%$ genome-wide ANI were used for database building ²¹. For fungal genomes,
779 all genomes downloaded from the MycoCosm project (on 24th Jan., 2022) were used ²².
780 The amino acid sequences of predicted gene on the genomes were obtained using
781 FragGeneScanRs. The Universal Single Copy Gene (USCG) gene set for GTDB
782 genomes were also extracted via hmmer software.

783 To test the performance of our pipeline, we specifically chose genomes that are
784 not included in the GTDB database (the database was used for graph building). In
785 particular, the bacterial/archaeal genomes, mostly MAGs, reported by Ye and colleagues
786 ²³ and Tara Ocean MAGs (total 8,466 MAGs) ²⁴ were used. We randomly selected 1000
787 genomes/MAGs from Ye's collection and use them as query genomes to test the
788 performance and accuracy of GSearch. To compare with other database search tools for
789 large database e.g. phage database, we compare GSearch with PhageCloud ²⁵, which

790 builds a graph database based on the labels of each phage genome (e.g., environment
791 source) and its search algorithm is Dashing2 (not published).

792

793 *Recall of AAI-, ANI- and MinHash-based nearest neighbor searching for*
794 *bacteria/archaea, fungi and bacteriophage genomes.*

795 To benchmark how GSearch performs compared to ANI/AAI- and MinHash-based tools,
796 we ran FastANI, Diamond blastp-based AAI and Mash to find the best neighbors for the
797 same query dataset and evaluate whether or not the best neighbors found by GSearch
798 were the same. FastANI parameters for the bacterial dataset were the following: fastANI
799 --ql query_path.txt --rl gtdb_path.txt -k 16 -p 24 --minFrac 3000 -o ANI.txt. GTDB
800 database was split into 50 subsets and run each subset parallelly on a multi-node
801 supercomputer to reduce memory requirement. MASH parameters were: mash sketch -a
802 (for AA only) -k 21 (7 for AA) -s 12000 -p 24 GTDB/*.fna > gtdb.msh; mash dist -p 24
803 gtdb.msh query.msh. For AAI calculation, the corresponding script in the enveomics
804 package ²⁶ was used: aai.rb -1 query.faa -2 db.faa -p diamond -t 24. Hmmer was used to
805 search for universal single copy gene against pre-built hmm profiles (120 for archaea and
806 122 for bacteria respectively); the profiles were obtained from the GTDB-tk software. For
807 bacteriophage genome, FastANI fragment size 1000 was used instead of 3000 while
808 aai.rb fragment size is 500 instead of 1000 with minimal number of matches of 5. MASH
809 kmer size 11 and 7 was used for nucleotides and amino acid, respectively, for
810 bacteriophage. For fungal genome ANI calculation, we use MUMMER v4.0.0 with default
811 parameters ²⁷. Gene prediction for fungal genomes was performed using GeneMark-ES
812 v2 (--fungus --ES) ²⁸. Kmer size 21 and 11 was used for fungal genomes in MASH for

813 nucleotides and amino acid, respectively. Detailed description of kmer size for each type
814 of genome can be found in Supplemental Note 2.

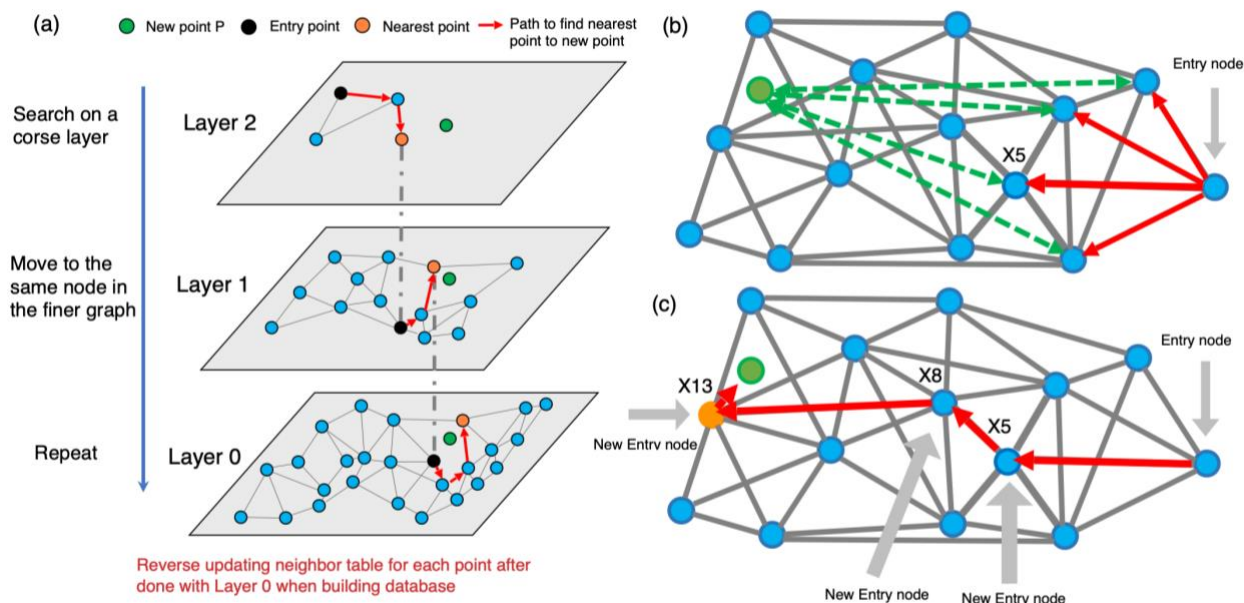
815 We calculated recall for our tool compare to standard ANI/AAI and MASH in the
816 following way: since biological species database are generally sparse because we are far
817 away from sequencing all species in the environment and likely the existence of natural
818 gaps in diversity, a larger top K by HNSW (e.g., 100) compared to the value used in
819 standard benchmark dataset will offer little, if any advantage, especially when the query
820 are relatively new, e.g. a new family compare to database genomes. Therefore, we use
821 top 5 and 10. Top 5 and top 10 recall are calculated based on top 5 and 10 neighbors
822 found by our tool and the available tools, and if all top 5 or 10 found by the latter tools
823 were also in top 5 or 10 of our tool, then recall was 100%. Similarly, if only 4 or 9 are
824 found by our tools, then recall was 80% and 90% respectively. However, if the distance
825 of query to some of the top 10 or top 5 neighbors found by our tool at the nucleotide level
826 was larger than 0.9850 for bacterial genomes, these matches will be filtered out and only
827 those neighbors below 0.9850 will be used (e.g. 8 out of 10 are kept, so only top 8 is
828 compared) because we have shown that above this threshold, Minhash-based methods
829 will lose accuracy and this is not specific to HNSW. Similar rules were applied for the
830 amino acid level searches with the threshold 0.9720 for filtering out bacterial genomes.

831 References

- 832 1. Ertl, O. ProbMinHash – A Class of Locality-Sensitive Hash Algorithms for the
833 (Probability) Jaccard Similarity. *IEEE Transactions on Knowledge and Data Engineering*,
834 1-1 (2020).
- 835 2. Broder, A.Z. in Proceedings. Compression and Complexity of SEQUENCES 1997 (Cat.
836 No. 97TB100171) 21-29 (IEEE, 1997).
- 837 3. Ondov, B.D. et al. Mash: fast genome and metagenome distance estimation using
838 MinHash. *Genome Biology* **17**, 132 (2016).
- 839 4. Koslicki, D. & Zabeti, H. Improving MinHash via the containment index with applications
840 to metagenomic analysis. *Applied Mathematics and Computation* **354**, 206-215 (2019).
- 841 5. Brown, C.T. & Irber, L. sourmash: a library for MinHash sketching of DNA. *Journal of*
842 *Open Source Software* **1**, 27 (2016).
- 843 6. Bovee, R. & Greenfield, N. Finch: a tool adding dynamic abundance filtering to genomic
844 MinHashing. *Journal of Open Source Software* **3**, 505 (2018).
- 845 7. Irber, L. et al. Lightweight compositional analysis of metagenomes with FracMinHash
846 and minimum metagenome covers. *bioRxiv*, 2022.2001.2011.475838 (2022).
- 847 8. Baker, D.N. & Langmead, B. Dashing: fast and accurate genomic distances with
848 HyperLogLog. *Genome Biology* **20**, 265 (2019).
- 849 9. Yang, D., Li, B., Rettig, L. & Cudré-Mauroux, P. D²histoSketch: Discriminative and
850 Dynamic Similarity-Preserving Sketching of Streaming Histograms. *IEEE Transactions*
851 *on Knowledge and Data Engineering* **31**, 1898-1911 (2019).
- 852 10. Rowe, W.P.M. et al. Streaming histogram sketching for rapid microbiome analytics.
853 *Microbiome* **7**, 40 (2019).
- 854 11. Ertl, O. Superminhash-A new minwise hashing algorithm for jaccard similarity estimation.
855 *arXiv preprint arXiv:1706.05698* (2017).
- 856 12. Ertl, O. in Proceedings of the 24th ACM SIGKDD International Conference on
857 Knowledge Discovery & Data Mining 1368–1377 (Association for Computing
858 Machinery, London, United Kingdom; 2018).
- 859 13. Ioffe, S. in 2010 IEEE International Conference on Data Mining 246-255 (2010).
- 860 14. Aumüller, M., Bernhardsson, E. & Faithfull, A. ANN-Benchmarks: A benchmarking tool
861 for approximate nearest neighbor algorithms. *Information Systems* **87**, 101374 (2020).
- 862 15. Déraspe, M., Boisvert, S., Laviolette, F., Roy, P.H. & Corbeil, J. Fast protein database as
863 a service with kAAmer. *bioRxiv*, 2020.2004.2001.019984 (2020).
- 864 16. Van der Jeugt, F., Dawyndt, P. & Mesuere, B. FragGeneScanRs: faster gene prediction
865 for short reads. *BMC Bioinformatics* **23**, 198 (2022).
- 866 17. Malkov, Y.A. & Yashunin, D.A. Efficient and Robust Approximate Nearest Neighbor
867 Search Using Hierarchical Navigable Small World Graphs. *IEEE Transactions on Pattern*
868 *Analysis and Machine Intelligence* **42**, 824-836 (2020).
- 869 18. Fu, C., Xiang, C., Wang, C. & Cai, D. Fast approximate nearest neighbor search with the
870 navigating spreading-out graph. *arXiv preprint arXiv:1707.00143* (2017).
- 871 19. Parks, D.H. et al. GTDB: an ongoing census of bacterial and archaeal diversity through a
872 phylogenetically consistent, rank normalized and complete genome-based taxonomy.
873 *Nucleic Acids Research* (2021).
- 874 20. Gregory, A.C. et al. Marine DNA viral macro-and microdiversity from pole to pole. *Cell*
875 **177**, 1109-1123. e1114 (2019).
- 876 21. Roux, S. et al. IMG/VR v3: an integrated ecological and evolutionary framework for
877 interrogating genomes of uncultivated viruses. *Nucleic acids research* **49**, D764-D775
878 (2021).

- 879 22. Grigoriev, I.V. et al. MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic*
880 *acids research* **42**, D699-D704 (2014).
- 881 23. Ye, L., Mei, R., Liu, W.-T., Ren, H. & Zhang, X.-X. Machine learning-aided analyses of
882 thousands of draft genomes reveal specific features of activated sludge processes.
883 *Microbiome* **8**, 1-13 (2020).
- 884 24. Nishimura, Y. & Yoshizawa, S. The OceanDNA MAG catalog contains over 50,000
885 prokaryotic genomes originated from various marine environments. *Scientific Data* **9**,
886 305 (2022).
- 887 25. Rangel-Pineros, G. et al. From Trees to Clouds: PhageClouds for Fast Comparison of~
888 640,000 Phage Genomic Sequences and Host-Centric Visualization Using Genomic
889 Network Graphs. *PHAGE* **2**, 194-203 (2021).
- 890 26. Rodriguez-R, L.M. & Konstantinidis, K.T. (PeerJ Preprints, 2016).
- 891 27. Marçais, G. et al. MUMmer4: A fast and versatile genome alignment system. *PLOS*
892 *Computational Biology* **14**, e1005944 (2018).
- 893 28. Ter-Hovhannisyanyan, V., Lomsadze, A., Chernoff, Y.O. & Borodovsky, M. Gene prediction
894 in novel fungal genomes using an ab initio algorithm with unsupervised training. *Genome*
895 *research* **18**, 1979-1990 (2008).
- 896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922

923 **Figures**



924 **Figure 1.** Schematic overview of GSearch building graph and searching graph steps.
925 (a) Graph was clasped into hierarchical layers following exponential decay probability.
926 In this graph, ef and M , represent the number of searches when finding nearest
927 neighbors and maximum allowed number of neighbors for each node, respectively. In
928 each layer, starting from an entry node (random or inherit from layer above it,
929 depending on whether it is the top layer or not), GSearch finds the closest connected
930 neighbor of the entry node and assigns it as the new entry point P (b), and then
931 traverses in a greedy manner (i.e., update the entry point using the newly found closest
932 connected neighbor (c)) until the nearest neighbor in the layer is found, and then goes
933 to next layer. This process is repeated until required number of nearest neighbors are
934 all found for the given new querying/inserting point. For building graph, after the
935 required number of nearest neighbors are found, a reverse update step will be
936 performed to update neighbor list of all nodes in the graph.

937
938
939
940
941
942

943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988

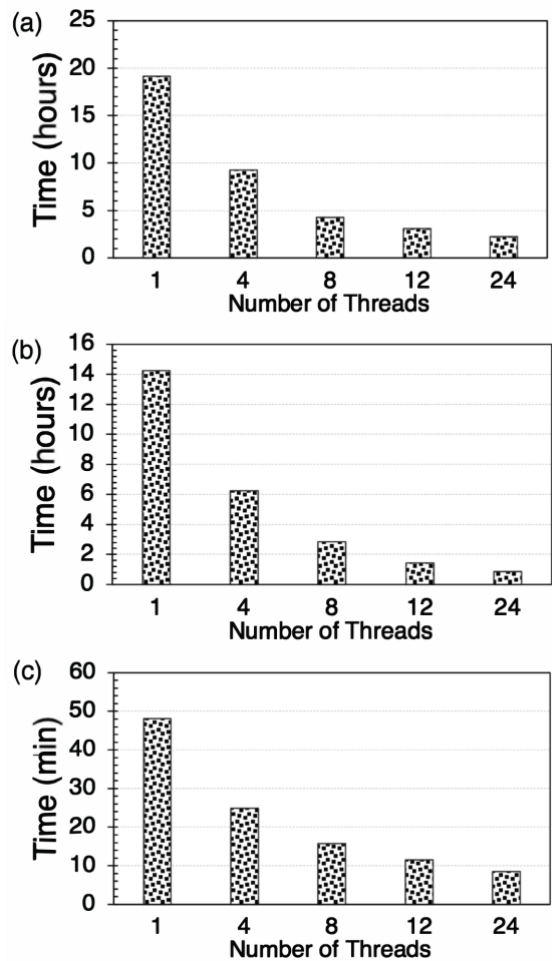
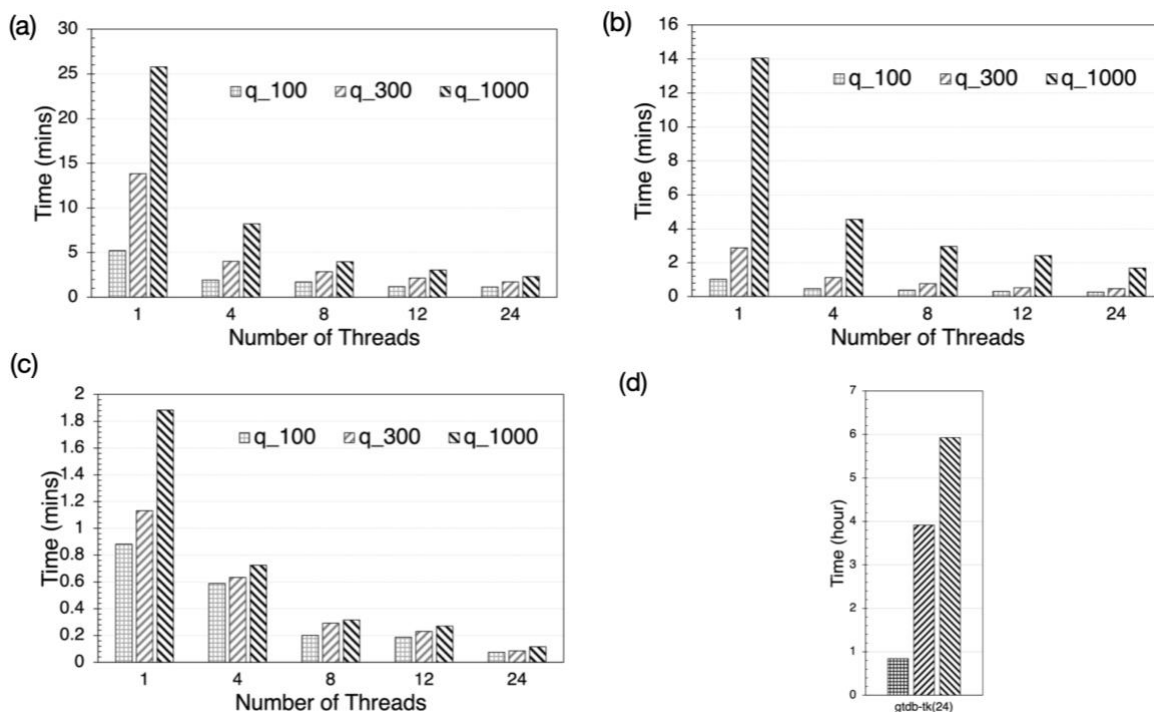


Figure 2. Scalability of database building process with the number of threads used. Panels show total wall time (y=axes) for building GTDB genome (nucleotide level) (a), whole-genome proteome (amino acid level) (b) and universal gene set proteome (c) databases. All tests were ran on a 24-thread Intel (R) Xeon (R) Gold 6226 processor, with 40GB memory available.



989

990

991

Figure 3. Total request time (wall time) for searching query genomes against the pre-built database of all GTDB genome (v207) at the whole-genome nucleotide (a), whole-genome proteome (b) and universal gene set proteome (c) levels. 100, 300 and 1000 query genomes (figure key) were used on a 24-thread Intel (R) Xeon (R) Gold 6226 processor. On average, database loading time ranged from 5-10 seconds. (d) is time needed to classify the same genomes using GTDB-Tk on the same 24-thread node.

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

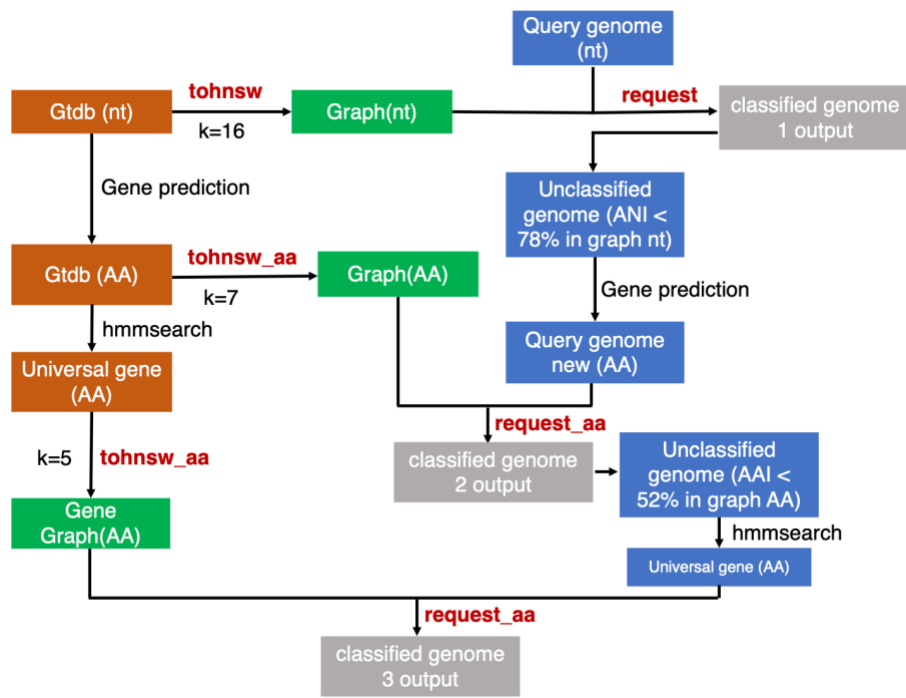
1012

1013

1014

1015

Figure 4. Overview of the GSearch pipeline for classifying prokaryotic genomes. Orange boxes denote steps that aim to prepare genome files, in different formats, for graph building while



1016 green boxes denote building steps of the graph database (in nucleotide or amino acid format).
 1017 Blue boxes indicate input/query genomes to search against the database while grey boxes
 1018 indicate classification output for each input. Gene prediction was done using FragGeneScanRs
 1019 and hmmsearch as part of the hmmer software for homology search. Two key steps of
 1020 GSearch: tohns_w (aa) and request (aa) are used to build graph database and request new
 1021 genomes, respectively. Two thresholds are used in the pipeline to decide between whole
 1022 nucleotide vs. whole-genome amino acid search and whole-genome amino acid vs. universal
 1023 gene amino acid, 78% ANI and 52% AAI, corresponding to Probinhash distance 0.9850 and
 1024 0.9375, respectively (see main text).

1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045

Tables

1046 **Table 1.** Request/search performance on major CPU platforms for GTDB v207 database for
1047 1000 queries.

CPU	Number of threads	Clock speed (GHz)	Request time for nt (min)	Gene Prediction-FGSrs (min) ^c	Request time for proteome (min)	hmmsearch time (min) ^d	Request time for USCG (min)
Intel (R) Xeon (R) Gold 6226 ^a	24	2.70	2.329	1.348	1.334	0.524	0.117
Intel (R) Core i7-7770HQ ^b	8	2.80	8.654	6.764	2.041	1.534	0.510
AMD EPYC 7513a ^a	32(24 used)	2.60	1.937	1.120	1.021	0.345	0.102
Apple M1 Pro ^b	10	3.22	2.369	2.12	0.866	0.498	0.168

1048
1049
1050
1051
1052
1053
1054
1055
1056

^a RHEL v7.9, Linux v3.10.0-1160, all threads used.

^b MacOS v12.3, Darwin 21.4.0, all threads used.

^c Parallel package was used to run multiprocessing at the same time. FGSrs stands for FragGeneScanRs. Note that in practice only those genomes failed in the Request for nt step (best found is less than 78% ANI) will be used in this step.

^d Only 100 genomes are used for testing hmmsearch because this step is for very new genomes at order level or above and we often do not have that many new genomes in a real-world dataset. Parallel Packages was used to run multiple processes of hmmsearch, one thread per process for hmmsearch.