

Optimized k -mer search across millions of bacterial genomes on laptops

Francesca Brunetti^{1,2,*}, Karel Břinda^{2,*}

¹ Department of Public Health and Infectious Diseases, Sapienza University of Rome, Rome, Italy

² Inria, Irisa, Univ. Rennes, Rennes, France

* francesca.brunetti@uniroma1.it, karel.brinda@inria.fr

ABSTRACT

Comprehensive bacterial collections have reached millions of genomes, opening new opportunities for point-of-care diagnostics and epidemiological surveillance. However, local real-time search over such collections on commodity hardware remains difficult. Currently, only LexicMap and Phylign enable local search and alignment at such a scale; among them, only Phylign is designed to run on laptops, via a subindex approach informed by phylogenetic compression. However, Phylign's performance deteriorates on long and divergent queries because it uses COBS as a k -mer-based prefilter before alignment with Minimap2. Meanwhile, recent k -mer indexes such as Fulgor and Themisto have emerged, but there is no practical methodology for selecting, combining, and parameterizing them for phylogenetically partitioned million-genome search under constraints.

Here, we develop an end-to-end methodology for k -mer matching in phylogenetically compressed bacterial collections. We formalize a matching strategy defined by matching mode, query type, and reference characteristics, and use this to shortlist candidate indexes and benchmark them under space–time trade-offs. As a case study, we address plasmid search over AllTheBacteria, compare multiple index types, and identify configurations optimizing the Pareto frontier of space and speed. Guided by these results, we implement a phylogenetically compressed variant of Fulgor, integrate it into Phylign, and obtain Phylign-Fulgor, a laptop-ready pipeline for million-genome search. On the 661k collection, Phylign-Fulgor makes the prefiltering step $\sim 4\times$ faster than Phylign-COBS at the cost of a $1.2\times$ larger index. On AllTheBacteria, its k -mer filter is $20\times$ – $300\times$ faster in real time than LexicMap's alignment-based search and uses $\sim 20\times$ smaller disk space. The full Phylign-Fulgor workflow including Minimap2 alignments is slower than LexicMap for a single plasmid but competitive or faster for batched plasmid queries. Phylign-Fulgor has comparable matching sensitivity to LexicMap, is less sensitive at the alignment level, but always stays within a laptop RAM budget ($\sim 5\times$ – $20\times$ lower memory than LexicMap).

KEYWORDS

k -mer indexing; phylogenetic compression; colored de Bruijn graphs; subindex search; million-genome search; Pareto optimization; plasmid surveillance

40 INTRODUCTION

41 Advances in DNA sequencing technologies have driven an exponential expansion of public microbial genome
42 archives [1–5]. The European Nucleotide Archive (ENA) alone grew from 661,405 curated bacterial genomes of
43 the “661k” Illumina snapshot [3] as of 2018 to over 2.4 million by August 2024, as captured in
44 AllTheBacteria [4]. Modern collections of assembled bacterial genomes such as AllTheBacteria [4], NCBI’s
45 bacterial assemblies, and GTDB [12] now contain thousands of species and millions of genomes. These
46 databases open up qualitatively new possibilities for point-of-care applications, including rapid AMR
47 diagnostics [8], plasmid and mobile-element tracing [9, 10], and large-scale epidemiological surveillance [11,
48 12]. Yet such applications ultimately rely on fast and sensitive local search and alignments across these million-
49 genome collections on commodity hardware.

50 Traditional alignment tools such as BLAST [13] remain highly sensitive but become impractical at this
51 scale: they typically require substantial memory and compute, which makes real-time search over millions of
52 genomes infeasible. A common alternative is to use k -mer indexes [14], which reduce genomes to k -mer sets
53 and search based on determining which query k -mers belong to which reference sets. However, even with k -mer
54 indexes, real-time local search across million-genome collections remains challenging on laptops. State-of-the-
55 art k -mer indexes either do not scale to millions of genomes [15, 16] or assume access to powerful clusters or
56 servers [17, 18]. Other large-scale search frameworks, such as BIGSI-like systems including BIGSI [17],
57 Metagraph [19], Logan Search [20], or general tools such as MMseqs2 [21] and RopeBWT3 [22], also rely on
58 substantial computational infrastructure. LexicMap provides fast alignment against multi-million-genome
59 bacterial collections but targets server-class hardware and multi-terabyte disk space [22].

60 To our knowledge, only LexicMap [23] and Phylign [5] currently provide sensitive local search and
61 alignment at the scale of million assembled genomes. LexicMap enables fast alignment across AllTheBacteria
62 and similar-scale collections but uses multi-terabyte indexes and large memory, which makes it unsuitable for
63 portable devices. Phylign, in contrast, is designed for laptops: it uses phylogenetic compression [5] to partition
64 genomes into batches of closely related, phylogenetically ordered genomes. Phylign builds separate COBS [23]
65 k -mer indexes per batch and searches them independently before aggregating results, optionally followed by
66 targeted Minimap2 [24] alignment for best-scoring matches. This subindex paradigm enabled, for the first time,
67 gene- and plasmid-level searches against the 661k collection [3] in hours on a laptop.

68 Over the last few years, more advanced k -mer indexes such as Fulgor [15], Themisto [16], and
69 Metagraph [27] have emerged. These tools offer substantially better scaling than COBS and in principle could
70 be used as phylogenetic subindexes to accelerate search on larger collections and improve sensitivity on long or
71 divergent queries. In practice, however, integrating them into subindex-based workflows is non-trivial. It
72 requires a principled way to map a biological question to a concrete k -mer matching strategy: which k -mers
73 must be matched, which matching modes are acceptable, and which index features (for instance, thresholding
74 schemes) are admissible under the subindex paradigm.

75 Despite the broad ecosystem of k -mer indexes now available [15–17, 23, 25–40], there is still no practical
76 end-to-end framework that (1) maps concrete biological questions to implementable k -mer matching strategies,
77 (2) selects compatible indexes and configurations under (laptop) hardware constraints, and (3) evaluates disk-
78 memory–time trade-offs when large collections must be partitioned into multiple subindexes, as in Phylign. This
79 gap – both methodological and implementation-related – makes it difficult to move beyond tool-by-tool
80 benchmarks and to design laptop-centric workflows for millions of genomes in realistic point-of-care use cases,
81 such as plasmid surveillance.

82 Here, we make four contributions. First, we introduce a formal framework for defining k -mer matching
83 strategies under the subindex paradigm, based on three elements – matching mode, query type, and reference
84 characteristics – which together determine admissible k -mer indexes and configurations. Second, we develop a
85 procedure based on Pareto-optimization for selecting k -mer indexes across phylogenetic batches, jointly
86 optimizing disk usage and search time at the collection scale. Third, we apply this framework to plasmid search
87 across AllTheBacteria, benchmarking COBS against three modern k -mer indexes – Fulgor, Themisto, and
88 Metagraph – and identify Fulgor as the best trade-off index. Fourth, we integrate a phylogenetically compressed
89 variant of Fulgor into Phylign, obtaining Phylign-Fulgor and evaluating it on the 661k and AllTheBacteria
90 collections under realistic laptop and server scenarios.

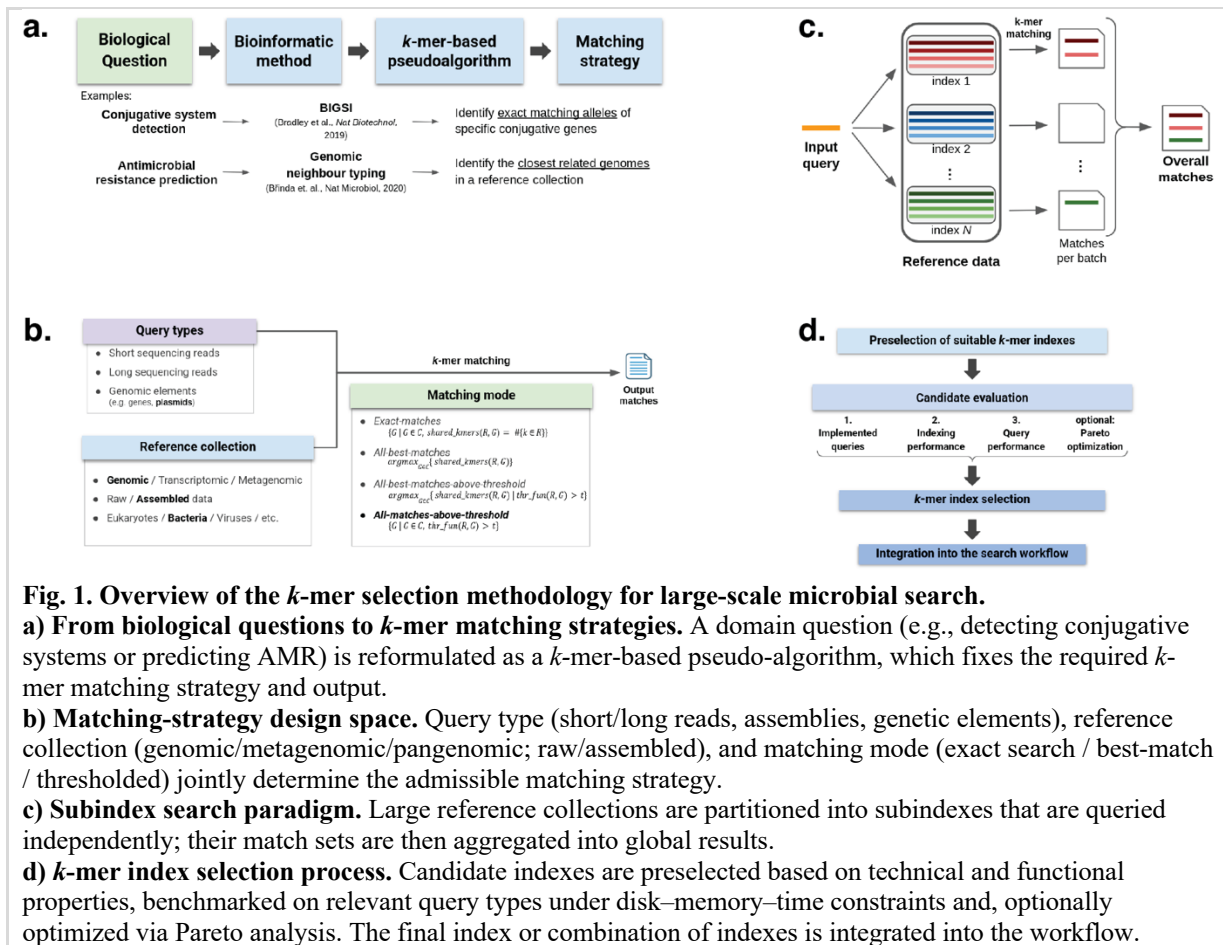


Fig. 1. Overview of the k -mer selection methodology for large-scale microbial search.

a) From biological questions to k -mer matching strategies. A domain question (e.g., detecting conjugative systems or predicting AMR) is reformulated as a k -mer-based pseudo-algorithm, which fixes the required k -mer matching strategy and output.

b) Matching-strategy design space. Query type (short/long reads, assemblies, genetic elements), reference collection (genomic/metagenomic/pangenomic; raw/assembled), and matching mode (exact search / best-match / thresholded) jointly determine the admissible matching strategy.

c) Subindex search paradigm. Large reference collections are partitioned into subindexes that are queried independently; their match sets are then aggregated into global results.

d) k -mer index selection process. Candidate indexes are preselected based on technical and functional properties, benchmarked on relevant query types under disk–memory–time constraints and, optionally optimized via Pareto analysis. The final index or combination of indexes is integrated into the workflow.

91 RESULTS

92 Methodology for optimizing local k -mer-based search

93 We developed a methodology for optimizing k -mer search on laptops (Fig. 1). Given a biological question
 94 (Fig. 1a), we assume it can be answered with k -mers and formalized as a pseudo-algorithm containing a k -mer
 95 matching step. Its formulation dictates a *matching strategy* defined by three components: *matching mode*,
 96 *reference collection characteristics*, and *query type* (Note S1).

97 Among the many available k -mer indexes [15–17, 23, 25–39], only a subset is relevant to any specific
 98 problem (Fig. 1b). Matching modes determine how k -mers are compared, query type (short/long reads, genes,
 99 plasmids, etc.) affects performance through length, quality, and divergence; and reference characteristics such as
 100 phylogenetic structure influence index size and search complexity – for example, clonal collections compress
 101 well and yield small, fast indexes, whereas heterogeneous ones inflate size and query time. Specific
 102 combinations of these elements define the matching strategy for a given biological question.

103 However, not every k -mer index can implement a chosen matching strategy under the *subindex paradigm*.
 104 To maintain searchability as collections grow faster than storage and memory, we rely on this paradigm as
 105 implemented in Phylign [5]: genomes are grouped into phylogenetically informed, highly compressible *batches*
 106 (Fig. 1c), each batch indexed and queried independently before aggregation. Complemented with phylogenetic
 107 structure [5], this enabled alignment across the full 661k collection on laptops [5]. This paradigm excludes some
 108 index designs. For instance, Themisto’s intersection regime evaluates thresholds only on k -mers present
 109 somewhere in the full index; k -mers absent in one batch but present in another break aggregation. The same
 110 limitation applies to Fulgor’s full-intersection method.

111 The final step is to identify the most suitable k -mer index from a shortlist of *candidates* via benchmarking
 112 (Fig. 1d, Note S2).

	COBS [23]	Fulgor [15, 25]	Themisto [16]	Metagraph [18]
a. Index properties				
exact indexing	NS	S	S	S
subindex compatibility	S	PS	PS	S
b. Matching modes				
exact-matches	S	NS	NS	S
best-matches	S	NS	NS	S
best-matches-above-threshold	S	NS	NS	S
all-matches-above-threshold	S	PS	PS	S

S: Supported; PS: Partially supported; NS: Not supported

Table 1. Comparison of functional capabilities of COBS, Fulgor, Themisto, and Metagraph.

Panels report support for exact indexing and subindex compatibility (a) and for individual matching modes (b). Functionality is classified as supported (S), partially supported or requiring minor code changes (PS), or not supported (NS). For details, see **Note S3**.

113 Case study: plasmid search over AllTheBacteria on laptops

114 Let us consider a specific use case: matching one or multiple plasmids on a laptop against the 661k/ATB
 115 collections. Plasmids are extrachromosomal DNA elements that facilitate horizontal gene transfer between
 116 bacterial isolates of the same or different species [42]. They can greatly vary in length and content, and their
 117 intrinsic heterogeneity undermines both k -mer-based filtering and alignment sensitivity at large scale. This
 118 mosaic architecture routinely breaks standard aligners – conserved regions inflate hit lists with non-informative
 119 matches, whereas divergent regions drop below sensitivity.

120 Here, we specifically consider the search of the EBI database of 2,826 plasmids, whose search was
 121 previously benchmarked in [5, 17], plus one additional plasmid from [22] used as a key benchmark dataset
 122 therein. Following previous work on plasmid matching using k -mers [5, 17], we adopt the *all-matches-above-*
 123 *threshold* matching mode between plasmid sequences and a reference collection of bacterial assemblies.

124 Technical and functional constraints on candidate k -mer indexes

125 We selected four indexes to benchmark for plasmid search within the subindex paradigm: COBS [23] as the
 126 current solution in Phylign, and Fulgor [15, 25], Themisto [16], and Metagraph [18] as potential replacements
 127 (**Tab. 1, Tab. S1**). To ensure sensitivity and accuracy, we considered only exact indexes.

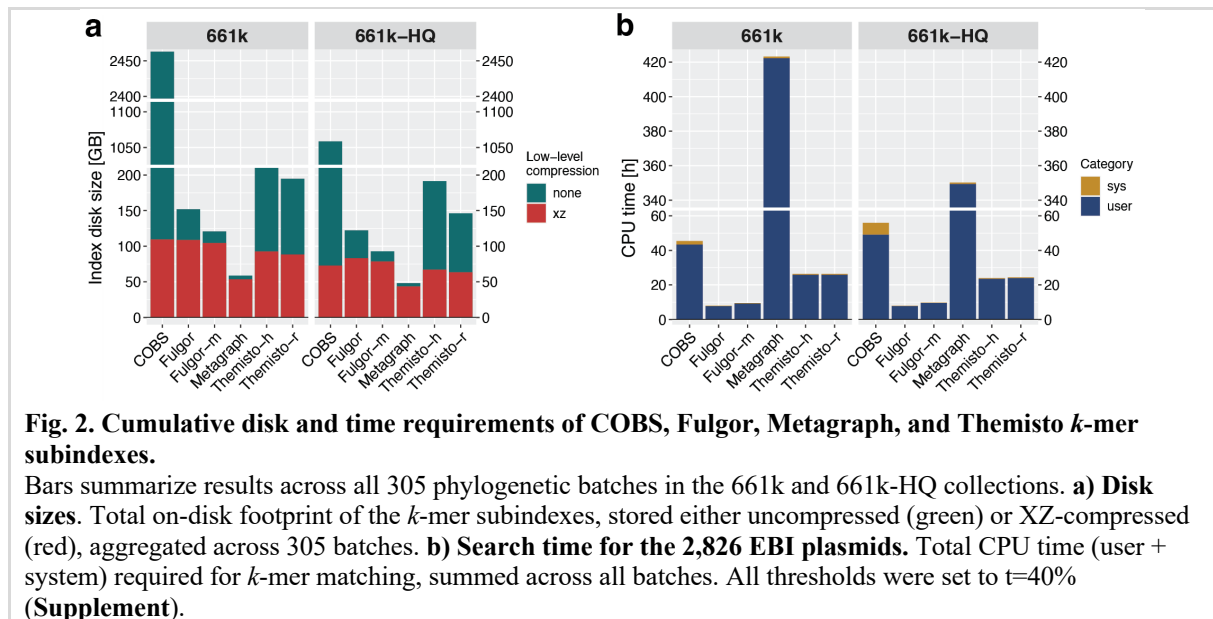
128 First, we assessed their compatibility with the subindex paradigm (**Tab. 1a**). Under this paradigm, the
 129 reference collection is partitioned into multiple batches, each batch is indexed and queried independently, and
 130 the batch-wise outputs are then aggregated into a global result.

131 Formally, let C be a reference collection, Q a query, and $M(C, Q)$ a well-defined k -mer-based similarity
 132 measure between C and Q that we wish to obtain. Let I_c denote an index built on the full collection C , and
 133 $\{I_j\}_{j=1\dots n}$ a family of subindexes forming a disjoint partition of C . Let $out_c(Q)$ be the output of searching Q
 134 against I_c , and $out_j(Q)$ the output against I_j . We say that an index family is *compatible with the subindex*
 135 *paradigm* if there exists an aggregation function $Comb$ such that, for all Q :

$$136 \quad out_c(Q) = Comb(out_{I_1}(Q), \dots, out_{I_n}(Q)) = M(C, Q)$$

137
 138 This rules out index designs whose semantics depend on the global k -mers presence rather than per subindex
 139 behavior.

141 Metagraph and COBS satisfy subindex compatibility directly through their threshold-based modes (termed
 142 *exact k -mer matching* and *approximate pattern matching* by their respective authors), returning all references
 143 whose similarity to the query exceeds a user-defined threshold computed over the full set of query k -mers.
 144 Themisto and Fulgor also implement thresholded search, but their default mode normalizes similarity positive k -
 145 mers only. This makes the denominator batch-dependent and breaks aggregation: each batch computes a
 146 different normalizing constant. To avoid this, both tools must operate in their full- k -mer regime and output
 147 explicit similarity values, which requires the “--include-unknown-kmers” parameter in Themisto and minor code
 148 changes in Fulgor. Under these conditions, the all-matches-above-threshold mode becomes implementable for
 149 all four tools (**Tab. 1b**).



150 *Disk–time trade-offs differ between index designs*

151 Next, we sought to understand the time-memory trade-offs of the four selected indexes when used within the
 152 subindex paradigm. For both Fulgor and Themisto, we considered their two color-annotation variants:
 153 Fulgor/meta-colored Fulgor (we will refer to it as Fulgor and Fulgor-m, respectively) and Themisto-
 154 hybrid/Themisto-roaring bitmap (Themisto-h and Themisto-r, respectively). For Metagraph, we considered only
 155 its row-diff relaxed BRWT variant, suggested by the authors as the best trade-off between compression and
 156 query performance. As the subindex paradigm enables storing subindexes in a compressed form and
 157 decompressing them on the fly, like implemented in Phylign with COBS, we primarily considered for each
 158 index both its variant compressed by XZ. Nevertheless, as k -mer indexes have substantially improved recently
 159 in internal compression [15, 16, 18], we also evaluated the index size without XZ compression.

160 As a reference database, we used the 661k collection, as extensive calibration and evaluation data are
 161 available in prior work [5] and 661k is already of sufficient size to be representative. We used the 305 genome
 162 batches that were inferred previously using MiniPhy [5], and embedded in the current version of Phylign
 163 (<https://github.com/karel-brinda/phylicn>). To account for the effects of data quality and contamination, we
 164 considered both its complete version (661k, $n=661,405$), as well as its subset omitting genomes that failed
 165 quality control (661k-HQ, $n=639,981$). In particular, the impact of omitting low-quality genomes has previously
 166 been shown as major for computational trade-offs [5]. We benchmarked the indexes on a single cluster node,
 167 with 30 CPUs and 200 GB of memory available.

168 First, we focused on disk sizes (Fig. 2a). The figure shows a consistent pattern across both collections.
 169 COBS produces by far the largest subindexes, even after XZ compression: its compressed footprint remains
 170 roughly an order of magnitude above all other tools. Fulgor and Fulgor-m form the smallest group overall, with
 171 XZ compression yielding only marginal additional savings. Themisto-h and Themisto-r occupy an intermediate
 172 band, with moderate reduction under XZ compression. Metagraph’s BRWT variant is extremely small even
 173 when uncompressed, suggesting that most structural redundancy has already been removed by the BRWT
 174 transform. The absolute values further highlight the divergence: across all 305 batches, COBS without
 175 compression is in the multi-terabyte range, Themisto around a few hundred gigabytes, and Fulgor consistently
 176 below that. Interestingly, all indexes except Metagraph have a similar size when they are compressed.

177 Second, we sought to understand differences in query performance at the level of individual indexes
 178 (Fig. 2b). Here, we observed substantial differences, and clear separation according to index types. The clear
 179 winner was Fulgor, which – both base and meta-colored – dominates with the lowest total CPU time across both
 180 661k and 661k-HQ. Themisto is consistently $\sim 2\times$ slower than Fulgor, regardless of variant. COBS is another
 181 $\sim 2\times$ slower than Themisto, confirming that its bit-sliced structure is penalized under heavy k -mer enumeration.
 182 Metagraph is an extreme outlier: even its optimized BRWT configuration is more than an order of magnitude
 183 slower than the others, reflecting the cost of decompression and traversal of a deeply compressed structure.

184 *Batch-level heterogeneity in space–time trade-offs*

185 To better understand index behavior beyond global averages, we examined performance on the 305
 186 phylogenetic batches of the 661k-HQ dataset (Fig. 3). A major advantage of the subindex paradigm is that

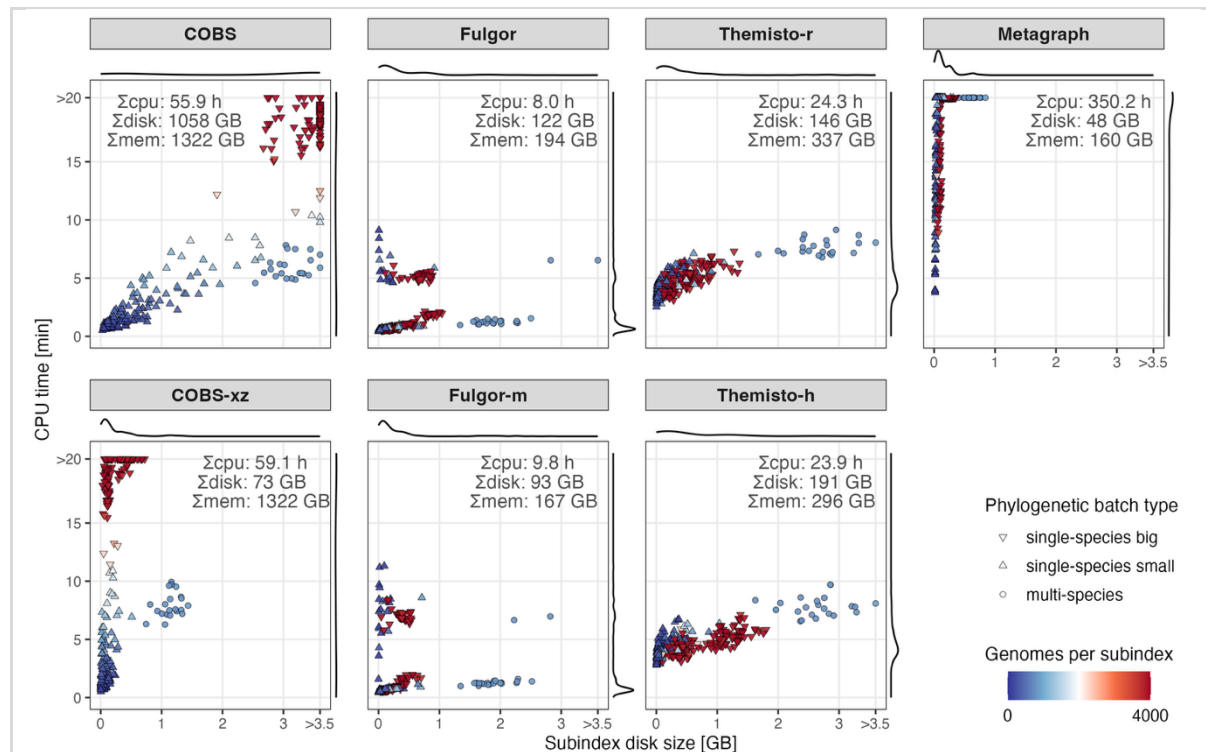


Fig. 3. Batch-level performance of k -mer subindexes in the 661k-HQ plasmid search.

For each index type, each point represents one phylogenetic batch; the x-axis shows the on-disk subindex size (GB), and the y-axis shows CPU time (min) to query the EBI plasmid set (same experiment as in Fig. 2).

Point color encodes the number of genomes per subindex, and point shape distinguishes large single-species, small single-species, and multi-species (“dustbin”) batches (threshold: 2,000 genomes). Text in each panel reports the total CPU time, total disk footprint, and total memory footprint across all batches. Marginal densities summarize the distributions of CPU time and disk size.

187 contributions of individual subindexes can be perfectly isolated and studied in isolation. We reused the data
 188 from the previous experiment, but slightly readjusted the formulation. As all indexes except COBS were
 189 sufficiently compressed on their own, we proceeded with seven subindex configurations: Fulgor, Fulgor-m,
 190 Themisto-r, Themisto-h, Metagraph, COBS, and COBS-xz (Supplement). In the latter (COBS-xz), we assume
 191 that the index is stored on disk in XZ-compressed form and decompressed on the fly, as implemented in
 192 Phylign.

193 The 661k-HQ collection contains 639,981 high-quality genomes from 2,336 species, with just 20 species
 194 accounting for 90% of genomes [3]. After phylogenetic compression, the top ten species – nearly 80% of the
 195 genomic content – occupy less than half of the database space, whereas “dustbin” batches of sparsely sampled
 196 species expand to proportions 9.4× larger than their precompression share [3]. To reflect this structure, MiniPhy
 197 defines 305 batches for 661k-HQ, of which 283 are single-species and 22 are multispecies dustbin batches [3].
 198 In Fig. 3, we further split single-species batches into large ($\geq 2,000$ genomes) and small ones, and examined
 199 disk–time trade-offs separately for these three batch classes.

200 We found that the index performance is strongly batch-dependent (Fig. 3). For all indexes, we see batches
 201 of the same type clustered. Large, low-diversity single-species batches tend to be best served by Fulgor/Fulgor-
 202 m, which are both small and fast, whereas high-diversity dustbin and small-species batches often favor COBS-
 203 xz, which remains time-competitive at modest size thanks to phylogenetic compression. Metagraph provides by
 204 far the best compression, but at the price of uncompetitive query times. Consistent with recent k -mer matching
 205 literature [42], all tools use substantially more memory compared to disk index size. No single index dominates
 206 across all phylogenetic batches, which motivates mixed-index optimization.

207 Pareto optimization of mixed index configurations

208 The large differences across index types prompted our interest in whether combining multiple index types
 209 within the same workflow provides measurable benefit. For instance, Fig. 3 suggests that high-diversity batches
 210 might benefit from COBS, whereas low-diversity ones could be served more efficiently with Fulgor. This can be
 211 formulated mathematically as follows: for a given disk space budget, how should we assign index types to
 212 individual batches so that they fit within the budget and simultaneously minimize the search time? This yields a
 213 standard multi-objective combinatorial optimization problem in the Pareto sense.

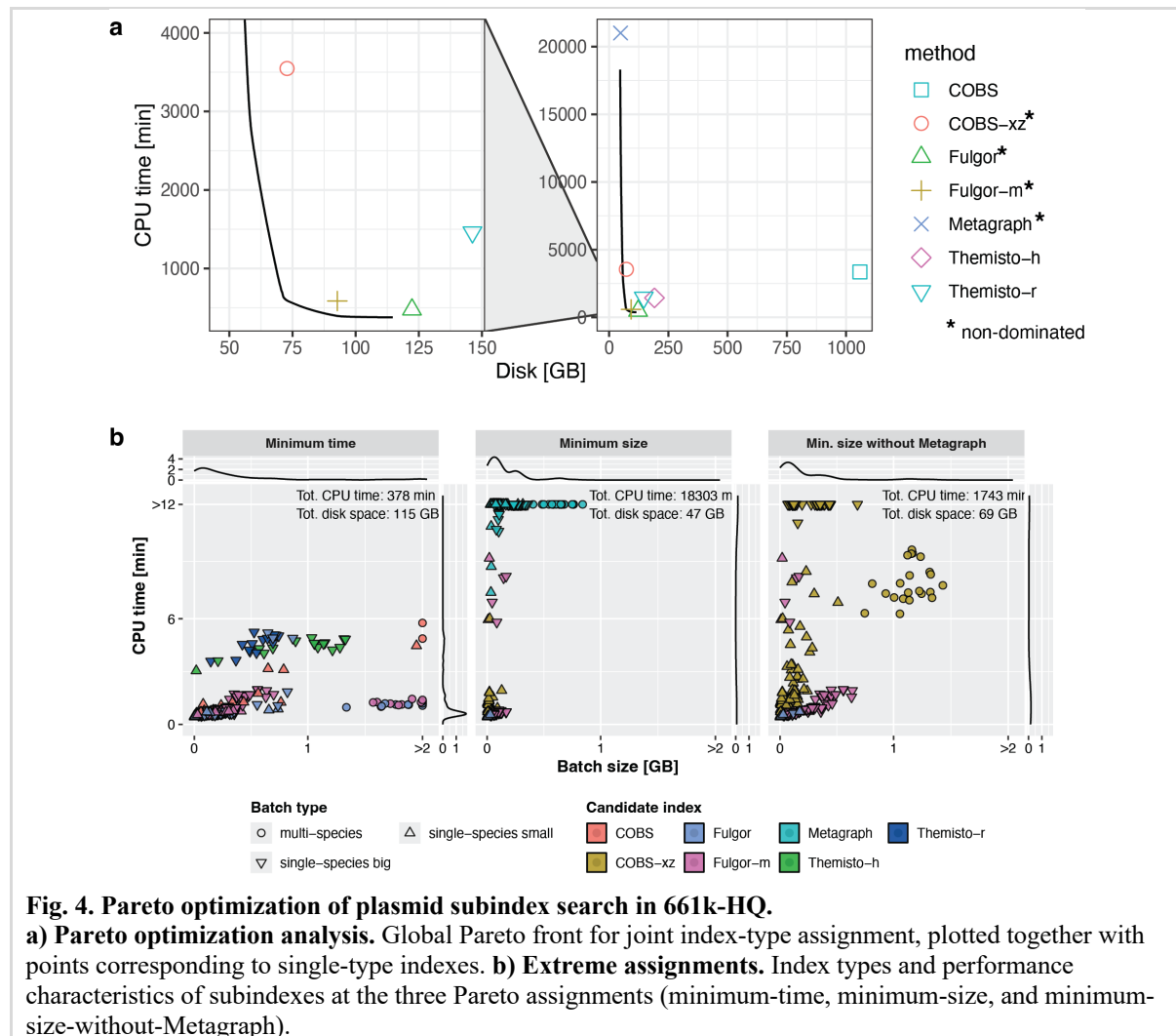


Fig. 4. Pareto optimization of plasmid subindex search in 661k-HQ.

a) Pareto optimization analysis. Global Pareto front for joint index-type assignment, plotted together with points corresponding to single-type indexes. **b) Extreme assignments.** Index types and performance characteristics of subindexes at the three Pareto assignments (minimum-time, minimum-size, and minimum-size-without-Metagraph).

214 Using the batch-level disk-time measurements, we computed the global Pareto front (Fig. 4a). For each
 215 batch and index type, we considered a cost vector (disk size, CPU time). First, we created two-objective Pareto
 216 sets for each batch, removing batch-level dominated configurations. Then, we computed the exact global Pareto
 217 front over all assignments by iteratively combining batches: starting from the zero-cost point (0 disk, 0 time), we
 218 repeatedly took the Minkowski sum of the current front, with the local Pareto set of the next batch and pruned
 219 globally dominated points. This yields the non-dominated combinations of index types over all batches under
 220 additive disk and time costs, in the standard setting of multi-objective combinatorial optimization [43].
 221 Conceptually, this procedure performs repeated Pareto sums of the batch-level Pareto sets, followed by
 222 dominance filtering, as studied for unions and Minkowski sums of non-dominated sets [44]. From these data, we
 223 also extracted three extreme assignments (Fig. 4b): the global minimum-time solution (fastest index for each
 224 batch), the global minimum-size solution (smallest index for each batch), and the minimum-size solution
 225 omitting Metagraph.

226 The examination of the front shows a characteristic steep descent: aggressive space minimization rapidly
 227 incurs extreme time penalties. The practical L-shaped region of the curve is dominated by three closely clustered
 228 candidates – COBS-xz, Fulgor-xz, and Fulgor – all lying near the knee of the efficient front and offering the
 229 strongest real-world trade-offs.

230 We then examined the index composition of the Pareto front endpoints (Fig. 4b). The fastest achievable
 231 solution required approximately 6 h and consumed 115 GB of disk space (Fig. 4b, left panel). Conversely,
 232 minimizing disk space reduces the requirements to 47 GB, but at a prohibitive cost of approximately 304 h of
 233 search time (Fig. 4b, central panel). Because these extremely slow configurations were consistently driven by
 234 Metagraph, we excluded all solutions containing Metagraph from further consideration.

235 Recomputing the Pareto front without Metagraph produced a more practical landscape. Disk requirements
 236 for indexing can be reduced to 69 GB of disk space (Fig. 4b, right panel), requiring a total of approximately
 237 26 h of search time, while the minimum search time remained unchanged relative to Fig. 4b (left panel). Among
 238 all single-type indexes, only COBS-xz and Fulgor-m fall between the two Pareto endpoints, making them the

239 strongest compromises between size and speed in practice: COBS is consistently smaller, whereas Fulgor is
240 consistently faster.

241 **Phylign-Fulgor and its experimental evaluation**

242 **Implementation and calibration**

243 Once subindex benchmarking and optimization are complete, the final step is to integrate the indexes into the
244 target workflow such as Phylign. This can be done either via a single index type, which is technically easier but
245 less efficient, or as a combination of multiple index types. As Phylign does not currently support a simultaneous
246 use of multiple index types, we focused on a single-index configuration. Based on our analyses, we selected
247 Fulgor-m (meta-colored Fulgor) as the base index. Due to Fulgor's limited native support for subindex search
248 (**Tab. 1**), we created its custom modification (a fork of v2.0.0), adapted to the subindex paradigm, with modified
249 parameters, a more suitable output format, and an improved control over k -mer matching (**Supplement**).

250 We developed Phylign-Fulgor, a fork of the original Phylign embedding Fulgor. This required a series of
251 source code modifications (**Supplement**). We included both the 661k-HQ ($n=639,981$ genomes, 305 batches)
252 and ATB-HQ collections (release v0.2, $n=1,858,610$ genomes, 650 batches). In both cases, phylogenetically
253 compressed Fulgor-m indexes were generated analogously to the previous benchmarks (Supplement). Phylign-
254 Fulgor is provided under the MIT license at <https://github.com/Francii-B/Phylign-Fulgor>, and all database files
255 are available on Zenodo (**Tab. S1**)

256 Phylign-Fulgor provides substantial dynamic range in its search sensitivity, primarily controlled by the k -
257 mer matching threshold parameter. To make its results comparable with LexicMap and Phylign-COBS as the
258 closest methods, we identified those thresholds that equalized for a single plasmid the number matched genomes
259 (**Fig. 5**). First, we used ATB-HQ to identify the maximum number of detectable matches by Phylign-Fulgor,
260 which was 498,410 corresponding to all sequences sharing at least one k -mer with the query ($t \approx 0.00005\%$).
261 Then, we used ATB-HQ collection to identify a threshold matching LexicMap's count (485,281 matches): this
262 was achieved via setting a threshold $t=0.007\%$, yielding 481,499 genomes (**Fig. 5**). Finally, we performed an
263 analogous calibration on the 661k-HQ collection
264 to reproduce the number of genomes matched by
265 Phylign-COBS (with its previously recommended
266 threshold $t=40\%$): this yielded 2,949 matches,
267 back-translated to Phylign-Fulgor threshold 18%
268 with 2,982 matches (**Fig. S1**).

269 **Experimental design**

270 We compared Phylign-Fulgor to the two existing
271 state-of-the-art tools locally capable of alignment
272 to million-genome collections: Phylign [5] and
273 LexicMap [22] (**Figs. 6, 7**). To evaluate the
274 impact of replacing COBS with Fulgor within
275 Phylign, we assessed performance under two
276 comparable experimental settings. First, we
277 reproduced on a standard laptop the plasmid-
278 search experiment described in the Phylign
279 paper [5], querying the EBI plasmids against the
280 661k-HQ collection. Second, we queried the
281 pCUVET18-1784 plasmid across the same 661k-
282 HQ dataset to compare the sensitivity of Phylign
283 and Phylign-Fulgor. Similarly, we evaluated the
284 performance of Phylign-Fulgor and LexicMap by
285 replicating the plasmid-search experiment from
286 the LexicMap paper [22], querying the
287 pCUVET18-1784 plasmid against the ATB-HQ
288 collection, and we additionally queried first 100
289 plasmid from the same EBI dataset to assess
290 performance across both laptop and server
291 environments..

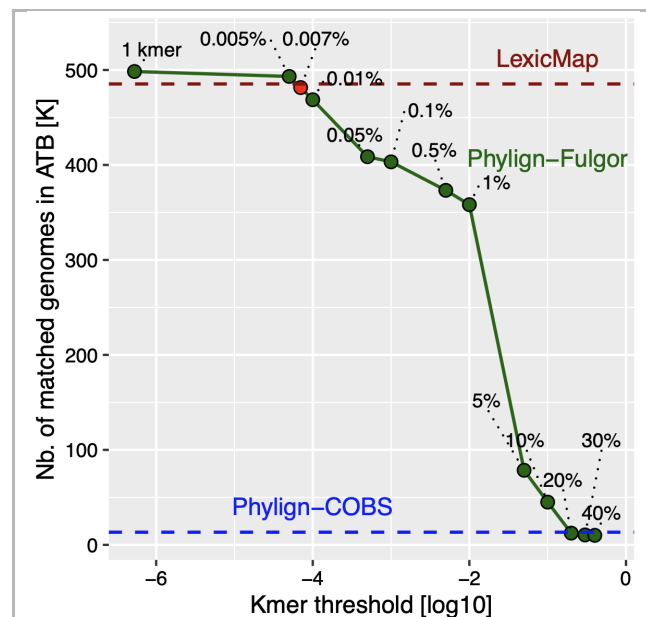
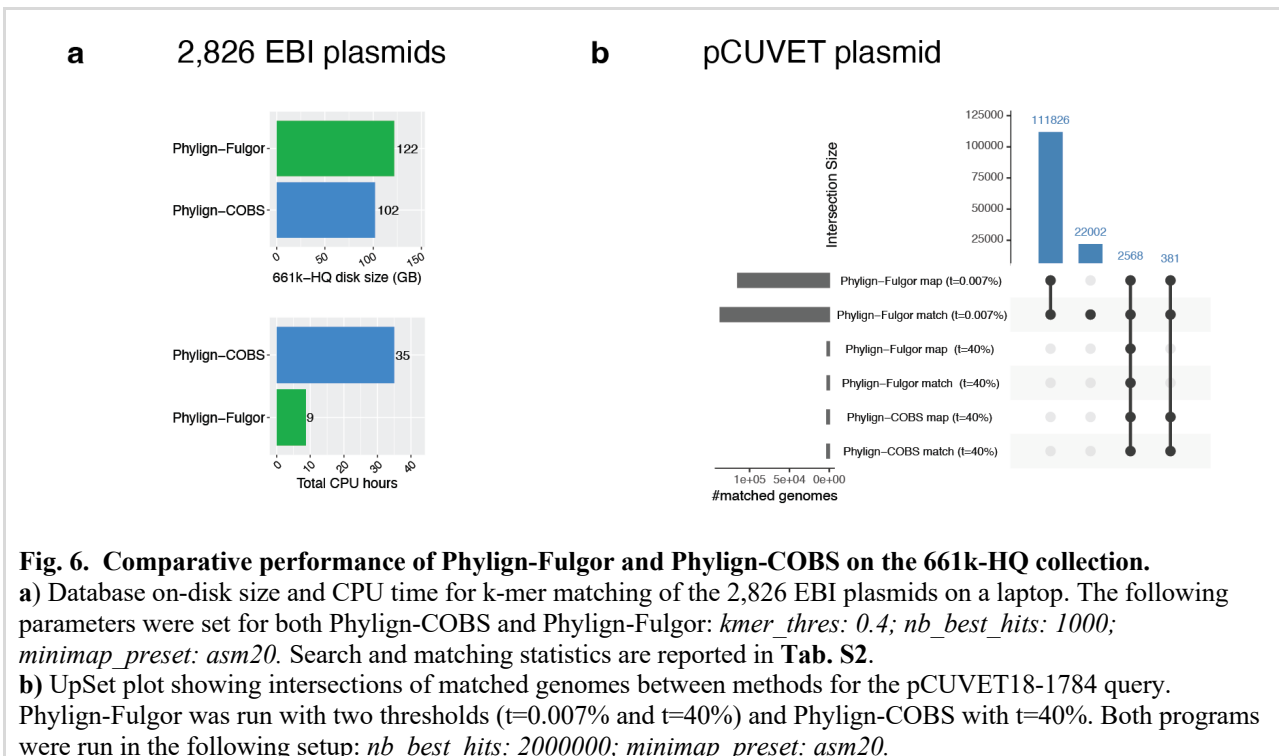


Fig. 5. Dynamic range of Phylign-Fulgor.

The graph shows the number of matched genomes for the pCUVET18-1784 plasmid in the ATB-HQ collection as a function of the minimum required proportion of matching k -mers in the *all-matches-above-threshold* mode, compared with LexicMap and Phylign-COBS (threshold $t=40\%$, as lower values cause a combinatorial explosion in COBS). The leftmost point corresponds to requiring a single matching k -mer; the red point ($t=0.007\%$) matches the LexicMap-calibrated threshold. The blue line reports Phylign-Fulgor results using across the ATB-HQ collection using the threshold calibrated relative to Phylign-COBS ($t=18\%$, $n=13,302$).



292 **Comparative performance of Phylign-Fulgor and Phylign-COBS**

293 First, we evaluated how the k -mer matching step in Phylign is affected by replacing COBS with Fulgor (**Fig. 6a**,
 294 **Tab. S2**). We used the k -mer matching experiments from the BIGSI benchmark [17], which was also the first
 295 step of the original Phylign benchmark [5], and queried the 2,826 EBI plasmids against the 661k-HQ collection
 296 on a laptop.

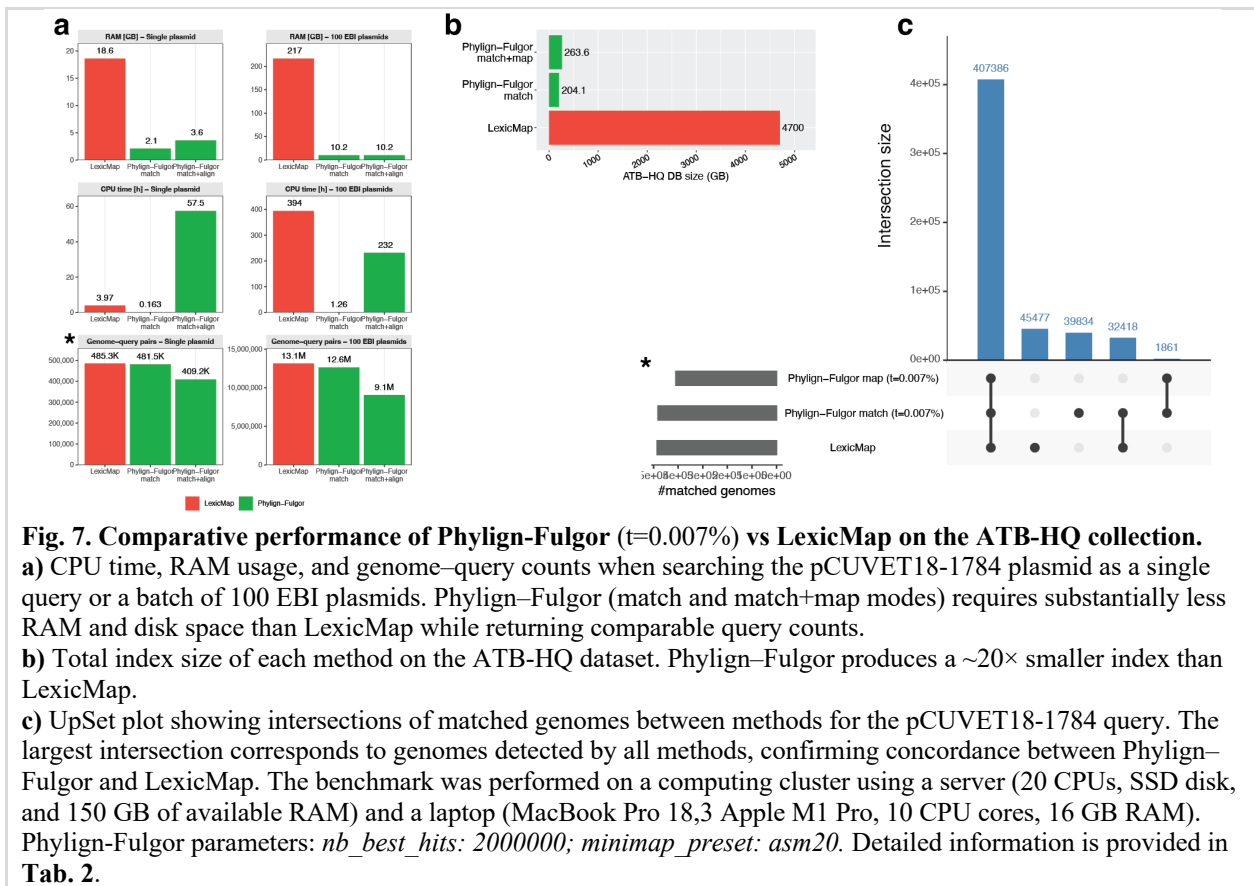
297 Phylign-Fulgor performed k -mer matching approximately $4\times$ faster than Phylign-COBS (8.4 h vs. 35 h of
 298 total CPU time), with a $1.2\times$ increase in its index size (from 102 GB to 122 GB) (**Fig. 6a**). However, because
 299 Fulgor implements exact k -mer indexing, it detected a factor of $2.5\times$ fewer matches than COBS (3,732,698 vs.
 300 9,163,123) (**Tab. S2**). This reduction did not substantially affect the runtime of the subsequent alignment step
 301 (Phylign-Fulgor: 14 CPU hours; Phylign-COBS: 15 CPU hours), but it did decrease the total number of
 302 alignments (7,705,870 vs. 8,980,408).

303 We next assessed how the COBS-to-Fulgor replacement affects sensitivity. For this analysis, we queried the
 304 pCUVET18-1784 plasmid across the 661k-HQ collection, reporting all detected matches (*nb_best_hits*:
 305 2000000), and compared the sets of matched genomes obtained before threshold calibration ($t=40\%$) and after
 306 calibration with respect to LexicMap ($t=0.007\%$).

307 Using the threshold previously recommended for Phylign-COBS ($t=40\%$), both tools identified a
 308 comparable number of matched genomes (2,568 for Phylign-Fulgor vs. 2,949 for Phylign-COBS), with all
 309 Phylign-Fulgor matches also found by Phylign-COBS (**Fig. 6b**). When the calibrated high-sensitivity threshold
 310 was used for Phylign-Fulgor ($t=0.007\%$), it detected 136,777 matches (approximately $46\times$ more than Phylign-
 311 COBS at $t=40\%$) while still recovering all matches detected by Phylign-COBS.

312 Despite the large variation in match counts across thresholds, all Phylign-Fulgor k -mer matching runs
 313 completed in 1–2 minutes of real time (4–6 minutes of total CPU time), making them 4–8 \times faster than Phylign-
 314 COBS in real time and 30–48 \times faster in total CPU time (Phylign-COBS real time: 10 minutes; total CPU time:
 315 ~2 h) (**Tab. S3**).

316 As the number of k -mer matches directly determines the workload of the subsequent alignment step, we
 317 also assessed how the increased match count affects alignment performance. With thresholds $t=40\%$, Phylign-
 318 Fulgor was $\sim 2\times$ faster than Phylign-COBS (~10 min vs. 18 min of real time; ~1.5 h vs 3.5 h of total CPU time),
 319 with only a modest reduction in detected alignments (22,956 vs 25,870) (**Tab. S3**).



320 Comparison to LexicMap on AllTheBacteria

321 As LexicMap is the current recommended tool for searching in AllTheBacteria, we sought to compare Phylign-
 322 Fulgor and LexicMap performance in a comparable setup. We evaluated two query types across the ATB-HQ
 323 collection: (1) the pCUVET18-1784 plasmid, replicating the plasmid-search experiment from the LexicMap
 324 paper [22], and (2) a batch of the first 100 plasmids from the same dataset. Because LexicMap is not scalable to
 325 portable devices, all LexicMap experiments were performed on a server node with an SSD disk and 20 CPUs
 326 available, whereas Phylign-Fulgor experiments were performed both on the server and additionally on a
 327 standard laptop, using its threshold calibrated with respect to LexicMap ($t=0.007\%$).

328 In contrast to the original Phylign (Phylign-COBS), where k -mer matching and alignment required similar
 329 amounts of time [5], Phylign-Fulgor’s k -mer matching step was hundreds of times faster than the subsequent
 330 Minimap2 alignment step ($\sim 830\times$ for the single plasmid and $\sim 180\times$ for the 100-plasmid batch). This is thanks to
 331 Fulgor’s high throughput and sensitivity, which generate many more matches in much less time.

332 When compared to LexicMap, Phylign-Fulgor k -mer matching was $\sim 3.1\times$ faster for the single plasmid and
 333 $\sim 28\times$ faster for the 100-plasmid batch, whereas the full Phylign-Fulgor workflow (k -mer matching + Minimap2
 334 alignment with the “asm20” preset) was $\sim 14\times$ slower for the single plasmid but $1.4\times$ faster for the 100-plasmid
 335 batch (**Fig. 7a**). This reflects the design of Phylign: it iterates over all phylogenetically compressed assemblies
 336 (unlike LexicMap, which keeps all assemblies uncompressed) and therefore has higher single-query costs, but
 337 these amortize once the batch is larger.

338 We then focused on Phylign-Fulgor space requirements in both memory and disk. On disk, its k -mer index
 339 required 204 GB, and the corresponding genome assemblies required ~ 59 GB, making its overall space
 340 requirements $17\text{--}23\times$ smaller than those of LexicMap (index size: 4.7 TB) (**Fig. 7b**). Phylign-Fulgor searches
 341 required substantially less memory than LexicMap, with reductions of $5\text{--}8\times$ for single-plasmid searches and $21\times$
 342 for the 100-plasmid batch (**Fig. 7a**).

343 Finally, we compared the sensitivity of both methods (**Fig. 7a,c; Tab. 2**). Phylign-Fulgor k -mer matching
 344 detected a comparable number of query–genome pairs to LexicMap for both the single plasmid (481,499 vs.
 345 485,281) and 100 plasmids (12,632,742 vs. 13,140,348) (**Fig. 7a**). After Minimap2 alignment, the number of
 346 query–genome pairs detected by Phylign-Fulgor decreased by $15\text{--}30\%$, resulting in $\sim 1.3\times$ fewer pairs than
 347 LexicMap for both the single plasmid ($n=409,247$) and the 100-plasmid batch ($n=9,069,409$). To assess
 348 sensitivity differences, we compared the set of genomes detected by Phylign-Fulgor and LexicMap for the

Method		Matching statistics			Laptop (10 cores)		Server (20 cores)		
		#target genomes	#genome-query pairs	#alignments	Real time [s]	CPU time [s]	Real time [s]	CPU time [s]	RAM [GB]
Single plasmid (pCUVET18-1784)									
Phylogn-Fulgor	match (1 k -mer)	498,410		–	102	205	193	552	2.5
	match ($t=0.007\%$)	481,499		–	104	209	259	585	2.1
	match+align ($t=0.007\%$)	409,247		1,635,746	14,841	114,534	11,372	206,856	3.6
LexicMap		485,281		7,192,512	(unscalable)		818	14,283	18.6
100 EBI plasmids									
Phylogn-Fulgor	match ($t=0.007\%$)	1,443,874	12,632,742	–	1,684	2,451	2,547	4,546	10.2
	match+align ($t=0.007\%$)	1,128,127	9,069,409	42,428,629	54,165	369,045	50,929	835,585	10.2
LexicMap		1,502,204	13,140,348	162,594,329	(unscalable)		71,900	1,418,475	216.6

Table 2. Detailed statistics for plasmid search in the ATB-HQ collection.

The configuration of the computers and search parameters are provided in the legend of **Fig. 7**.

349 single plasmid. Among the 481,499 matched genomes, Phylogn-Fulgor identified 90% ($n=439,804$) of those
 350 detected by LexicMap (**Fig. 7c**).

351 Overall, Phylogn-Fulgor trades a small loss in alignment-level sensitivity for dramatic improvements in
 352 space usage (20× smaller), memory (5–21× lower), and k -mer filtering throughput (3–28× faster), while
 353 becoming competitive or faster than LexicMap for realistic batches.

354 DISCUSSION AND CONCLUSIONS

355 Modern collections of millions of bacterial genomes may revolutionize point-of-care applications, yet the lack
 356 of local methods for alignment using portable devices is a critical gap. Here, we focused on k -mer-based search
 357 and developed a practical end-to-end methodology for combining modern k -mer indexes with phylogenetic
 358 compression [5] in local search via subindexes. We plugged this into Phylogn, obtaining the Phylogn-Fulgor
 359 pipeline, and showed that it can search such large collections on laptops within several minutes. As k -mer
 360 indexes have so far been mostly studied in isolation and applied on a large scale only on servers, this represents
 361 a major methodological advance.

362 One key observation is that index selection is not a single-metric choice, but a Pareto optimization problem
 363 balancing space, latency, batch behavior, and compatibility with the subindex paradigm. While this is intuitive,
 364 to the best of our knowledge, this is the first time to show this quantitatively. Across our plasmid-search use
 365 case, Fulgor, COBS, Themisto and Metagraph exhibited surprisingly different performance at the batch level
 366 depending on phylogenetic relatedness and batch size, reinforcing that cumulative high-level numbers for entire
 367 collections are not informative of index differences. Furthermore, most indexes are currently technically
 368 unsuited for the type of matching that would be needed, and source code modifications would almost always be
 369 necessary should they be used in practice.

370 Our study has several limitations. First, we relied on predefined MiniPhy batches, which fixes the
 371 phylogenetic granularity, which may not be optimal across all types of searches, taxonomic groups, or indexes.
 372 Given current compressive capabilities of Fulgor, one might imagine, for instance, building species-level
 373 indexes instead. Second, we used only the default parameter choice of each index, such as Fulgor’s minimizer
 374 size, yet these parameters may have a major effect on the resulting performance and could be easily included in
 375 the Pareto optimization. Third, we tested single type of queries in a constrained use case; yet other biological
 376 applications (e.g., metagenomic screening) may stress indexes in different ways. Finally, we did not evaluate the
 377 most recent Fulgor version, which substantially improved performance further, as this would require redoing our
 378 modification in its source code.

379 Future work will therefore focus on optimizing internal index parameters, co-optimizing batch construction
 380 with phylogenetic structure, and characterizing how species composition shapes size-speed trade-offs. We
 381 anticipate that the broader community may provide standard APIs in k -mer indexes, which would better allow us
 382 to integrate high-performance indexes in workflows such as ours without source code modifications. Ultimately,
 383 our goal is a fully adaptive version of Phylogn with support for databases, which could select, per task, the
 384 optimal combination of indexes and parameters to achieve near-optimal search, with the goal of providing near-
 385 real time search across millions of genomes on commodity hardware. Given the accelerating growth of bacterial
 386 genome databases and the demand for local data analyses, the methodological principles established here –
 387 pragmatic index evaluation and pareto-guided selection, combined with phylogenetic compression – provide a
 388 foundation for the next generation of large-scale k -mer search tools.

389 ACKNOWLEDGEMENTS

390 This research was supported by the French National Research Agency (ANR) under Grant ANR-24-CE45-1226
391 for the REALL project. Portions of this research were conducted at the GenOuest bioinformatics core facility
392 (<https://www.genouest.org>).

393 REFERENCES

- 394
395
396 1. Loh P-R, Baym M, Berger B (2012) Compressive genomics. *Nat Biotechnol* 30:627–630
- 397 2. Stephens ZD, Lee SY, Faghri F, Campbell RH, Zhai C, Efron MJ, Iyer R, Schatz MC, Sinha S, Robinson
398 GE (2015) Big Data: Astronomical or Genomical? *PLoS Biol* 13:e1002195
- 399 3. Blackwell GA, Hunt M, Malone KM, Lima L, Horesh G, Alako BTF, Thomson NR, Iqbal Z (2021)
400 Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences. *PLoS Biol*
401 19:e3001421
- 402 4. Hunt M, Lima L, Anderson D, Hawkey J, Shen W, Lees J, Iqbal Z (2024) AllTheBacteria - all bacterial
403 genomes assembled, available and searchable. *bioRxiv* 2024.03.08.584059
- 404 5. Brinda K, Lima L, Pignotti S, Quinones-Olvera N, Salikhov K, Chikhi R, Kucherov G, Iqbal Z, Baym M
405 (2025) Efficient and robust search of microbial genomes via phylogenetic compression. *Nat Methods*
406 22:692–697
- 407 6. Parks DH, Chuvochina M, Rinke C, Mussig AJ, Chaumeil P-A, Hugenholtz P (2022) GTDB: an ongoing
408 census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and
409 complete genome-based taxonomy. *Nucleic Acids Res* 50:D785–D794
- 410 7. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P (2020) A complete domain-
411 to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* 38:1079–1086
- 412 8. Brinda K, Callendrello A, Ma KC, MacFadden DR, Charalampous T, Lee RS, Cowley L, Wadsworth CB,
413 Grad YH, Kucherov G, O’Grady J, Baym M, Hanage WP (2020) Rapid inference of antibiotic resistance
414 and susceptibility by genomic neighbour typing. *Nat Microbiol* 5:455–464
- 415 9. Robertson J, Bessonov K, Schonfeld J, Nash JHE (2020) Universal whole-sequence-based plasmid typing
416 and its utility to prediction of host range and epidemiological surveillance. *Microb Genom* 6:.
417 <https://doi.org/10.1099/mgen.0.000435>
- 418 10. Partridge SR, Kwong SM, Firth N, Jensen SO (2018) Mobile genetic elements associated with
419 antimicrobial resistance. *Clin Microbiol Rev* 31:.
<https://doi.org/10.1128/CMR.00088-17>
- 420 11. Gardy JL, Loman NJ (2017) Towards a genomics-informed, real-time, global pathogen surveillance
421 system. *Nat Rev Genet*. <https://doi.org/10.1038/nrg.2017.88>
- 422 12. Baker KS, Jauneikaite E, Hopkins KL, Lo SW, Sánchez-Busó L, Getino M, Howden BP, Holt KE, Musila
423 LA, Hendriksen RS, Amoako DG, Aanensen DM, Okeke IN, Egyir B, Nunn JG, Midega JT, Feasey NA,
424 Peacock SJ, SEDRIC Genomics Surveillance Working Group (2023) Genomics for public health and
425 international surveillance of antimicrobial resistance. *Lancet Microbe* 4:e1047–e1055
- 426 13. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol*
427 *Biol* 215:403–410
- 428 14. Marchet C, Boucher C, Puglisi SJ, Medvedev P, Salson M, Chikhi R (2021) Data structures based on -mers
429 for querying large collections of sequencing data sets. *Genome Res* 31:1–12
- 430 15. Fan J, Khan J, Singh NP, Pibiri GE, Patro R (2024) Fulgor: a fast and compact *k*-mer index for large-scale
431 matching and color queries. *Algorithms Mol Biol* 19:3

- 432 16. Alanko JN, Vuhoitoniemi J, Mäklin T, Puglisi SJ (2023) Themisto: a scalable colored k-mer index for
433 sensitive pseudoalignment against hundreds of thousands of bacterial genomes. *Bioinformatics* 39:i260–
434 i269
- 435 17. Bradley P, den Bakker HC, Rocha EPC, McVean G, Iqbal Z (2019) Ultrafast search of all deposited
436 bacterial and viral genomic data. *Nat Biotechnol* 37:152–159
- 437 18. Karasikov M, Mustafa H, Danciu D, Kulkov O, Zimmermann M, Barber C, Rättsch G, Kahles A (2025)
438 Efficient and accurate search in petabase-scale sequence repositories. *Nature*.
439 <https://doi.org/10.1038/s41586-025-09603-w>
- 440 19. Chikhi R, Raffestin B, Korobeynikov A, Edgar R, Babaian A (2024) Logan: Planetary-scale genome
441 assembly surveys life’s diversity. *bioRxiv*
- 442 20. Steinegger M, Söding J (2017) MMseqs2 enables sensitive protein sequence searching for the analysis of
443 massive data sets. *Nat Biotechnol* 35:1026–1028
- 444 21. Li H (2024) BWT construction and search at the terabase scale. *Bioinformatics* 40:.
445 <https://doi.org/10.1093/bioinformatics/btae717>
- 446 22. Shen W, Lees JA, Iqbal Z (2025) Efficient sequence alignment against millions of prokaryotic genomes
447 with LexicMap. *Nat Biotechnol*. <https://doi.org/10.1038/s41587-025-02812-8>
- 448 23. Bingmann T, Bradley P, Gauger F, Iqbal Z (2019) COBS: A Compact Bit-Sliced Signature Index. In:
449 *String Processing and Information Retrieval*. Springer International Publishing, Cham, pp 285–303
- 450 24. Li H (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34:3094–3100
- 451 25. Pibiri GE, Fan J, Patro R (2024) Meta-colored Compacted de Bruijn Graphs. *Research in Computational*
452 *Molecular Biology* 131–146
- 453 26. Campanelli A, Pibiri GE, Patro R (2025) Fast pseudoalignment queries on compressed colored de Bruijn
454 graphs
- 455 27. Lemane T, Lezsoche N, Lecubin J, Pelletier E, Lescot M, Chikhi R, Peterlongo P (2024) Indexing and real-
456 time user-friendly queries in terabyte-sized complex genomic datasets with kindex and ORA. *Nature*
457 *Computational Science* 4:104–109
- 458 28. Marchet C, Limasset A (2023) Scalable sequence database search using partitioned aggregated Bloom
459 comb trees. *Bioinformatics* 39:i252–i259
- 460 29. Pandey P, Almodaresi F, Bender MA, Ferdman M, Johnson R, Patro R (2018) Mantis: A Fast, Small, and
461 Exact Large-Scale Sequence-Search Index. *Cell Systems* 1–7
- 462 30. Yu Y, Liu J, Liu X, Zhang Y, Magner E, Lehnert E, Qian C, Liu J (2018) SeqOthello: querying RNA-seq
463 experiments at scale. *Genome Biol* 19:167
- 464 31. Holley G, Melsted P (2020) Bifrost: highly parallel construction and indexing of colored and compacted de
465 Bruijn graphs. *Genome Biol* 21:249
- 466 32. Cracco A, Tomescu AI (2023) Extremely fast construction and querying of compacted and colored de
467 Bruijn graphs with GGCAT. *Genome Res* 33:1198–1207
- 468 33. Hernandez-Courbevoie Y, Salson M, Bessièrè C, Xue H, Gautheret D, Marchet C, Limasset A (2026)
469 REINDEER2: Practical abundance index at scale. In: *Lecture Notes in Computer Science*. Springer Nature
470 Switzerland, Cham, pp 156–171
- 471 34. Solomon B, Kingsford C (2016) Fast search of thousands of short-read sequencing experiments. *Nat*
472 *Biotechnol* 34:300–302
- 473 35. Seiler E, Mehringer S, Darvish M, Turc E, Reinert K (2021) Raptor: A fast and space-efficient pre-filter for
474 querying very large collections of nucleotide sequences. *iScience* 24:102782

- 475 36. Marchet C, Iqbal Z, Gautheret D, Salson M, Chikhi R (2020) REINDEER: efficient indexing of k -mer
476 presence and abundance in sequencing datasets. *Bioinformatics* 36:i177–i185
- 477 37. Harris RS, Medvedev P (2020) Improved representation of sequence bloom trees. *Bioinformatics* 36:721–
478 727
- 479 38. Bray NL, Pimentel H, Melsted P, Pachter L (2016) Near-optimal probabilistic RNA-seq quantification. *Nat*
480 *Biotechnol* 34:525–527
- 481 39. Almodaresi F, Sarkar H, Srivastava A, Patro R (2018) A space and time-efficient index for the compacted
482 colored de Bruijn graph. *Bioinformatics* 34:i169–i177
- 483 40. Campanelli A, Pibiri GE, Fan J, Patro R (2024) Where the patterns are: Repetition-aware compression for
484 colored de Bruijn graphs. *J Comput Biol* 31:1022–1044
- 485 41. Pitout JDD, Chen L (2023) The Significance of Epidemic Plasmids in the Success of Multidrug-Resistant
486 Drug Pandemic Extraintestinal Pathogenic *Escherichia coli*. *Infectious Diseases and Therapy* 12:1029–
487 1041
- 488 42. Sladký O, Veselý P, Břinda K (2025) FroM Superstring to Indexing: a space-efficient index for
489 unconstrained k -mer sets using the Masked Burrows-Wheeler Transform (MBWT). *Bioinform Adv*.
490 <https://doi.org/10.1093/bioadv/vbaf290>
- 491 43. Ehrgott M, Gandibleux X (2000) A survey and annotated bibliography of multiobjective combinatorial
492 optimization. *OR Spectr* 22:425–460
- 493 44. Klamroth K, Lang B, Stiglmayr M (2024) Efficient dominance filtering for unions and Minkowski sums of
494 non-dominated sets. *Comput Oper Res* 163:106506
- 495 45. Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, Forster J, Lee S, Twardziok
496 SO, Kanitz A, Wilm A, Holtgrewe M, Rahmann S, Nahnsen S, Köster J (2021) Sustainable data analysis
497 with Snakemake. *F1000Res* 10:33
- 498

499 Supplement

500

501 METHODS

502 Data and program preparation

503 **661k collection**

504 All genomes from the 661k collection (n=661,309) [3] were downloaded from Zenodo using Phylign (command
505 “*make download_asms*”) [5], as 305 phylogenetically compressed .tar.xz archives containing FASTA files. The
506 corresponding Zenodo records are listed in **Tab. S1**. Phylogenetically ordered lists of genomes were generated
507 for each batch and stored as text files using the following command:

508

```
509 tar -tf {batch}.tar.xz > {batch_related_genome_list}.txt
```

510

511 HQ genomes were extracted from the 661k collection [3], based on its original metadata and analyzed as a
512 separate collection. The HQ information was used as reported in *File4_QC_characterisation_661K.txt*
513 (**Tab. S1**). The batches and phylogenetic orderings were used as defined in [5].

514 **AllTheBacteria**

515 Genomes of the ATB collection from release 0.2 (n=2,440,377) [4] were downloaded as 665 phylogenetically
516 compressed .tar.xz archives of FASTA files from OSF (**Tab. S1**), together with their metadata stored in
517 *file_list.all.20240805.tsv.gz*. Phylogenetically ordered lists of all genomes, as well as for the HQ subsets
518 (n=1,858,610), were generated for each batch, following the same procedure described for the 661k collection.
519 The HQ information was obtained from the file *hq_set.sample_list.txt.gz*, available on OSF.

520 **Used indexes**

521 We used four *k*-mer indexes selected for evaluations:

- 522 • Fulgor v. 2.0.0 [15, 25], the default and the meta-colored variants (Fulgor and Fulgor-m,
523 respectively).
- 524 • Themisto v.3.2.2 [16], both hybrid and roaring bitmap variants (Themisto-h and Themisto-r,
525 respectively).
- 526 • Metagraph v.0.3.6 [7] (row-diff relaxed Multi-BRWT variant).
- 527 • COBS v.0.3.1 [23].

528 All indexes were downloaded from their respective GitHub repositories (**Tab. S1**).

529 **Evaluation of index space requirements**

530 Evaluations of index space requirements were performed on the GenOuest cluster (<https://www.genouest.org>),
531 using a single node with 30 CPUs and 200 GB of memory available for computation.

532 For each considered index (except COBS, which was already evaluated in [5]), we developed a dedicated
533 Snakemake [9] workflow to generate batch-level compressed indexes for each batch of the 661k and 661k-HQ
534 collections (305 batches in both cases). All indexes were created such that the phylogenetic reordering is used at
535 its max: the batches were phylogenetic and we retained the phylogenetic order in the input for each program.
536 The specific commands used for generating each subindex are reported below.

537 **Fulgor and Fulgor-m:**

538 For each batch, we built both Fulgor index types [15, 25] via a single invocation of the following command:

```
539 fulgor build -k 31 -m 20 -l {batch_genome_list} -o {fulgor_index_name} --check --meta
```

540

541 **Themisto-h and Themisto-r:**

542 We built the Themisto [16] indexes using the following commands:

543

```
544 themisto build -k 31 -d 20 -i {batch_genome_list} -o {themisto_index_name} -s {annotation_type}
```

545

546 where *{annotation_type}* indicated the specific type of data structure to use for the color annotation, namely
547 “roaring” to build Themisto-roaring bitmap indexes (Themisto-r) or “sdsl-hybrid” to generate Themisto-hybrid
548 indexes (Themisto-h).
549

550 **Metagraph**

551 We considered the Row-diff relaxed Multi-BRWT variant of Metagraph [7]. The following commands were used
552 to generate the DBG, its annotation matrix and to convert this last into row-diff relaxed Multi-BRWT format:

```
553  
554 cat {batch_genome_list} | metagraph build -k 31 --graph succinct --state small -o  
555 {metagraph_index_name}.dbg
```

556
557 where *{metagraph_index_name}.dbg* corresponds to the output DBG.
558

```
559 cat {batch_genome_list} | metagraph annotate -i {metagraph_index_name}.dbg --anno-filename -o  
560 {metagraph_column_annotation}
```

561
562 where *{metagraph_index_annotation}* corresponds to the file name of the output annotation matrix.
563

```
564 find -name "{metagraph_column_annotation}" | metagraph transform anno --anno-type row_diff --  
565 row-diff-stage {n} -i {metagraph_index_name}.dbg -o {metagraph_rowdiff_annotation}
```

```
566  
567 find -name "{metagraph_rowdiff_annotation}" | metagraph transform anno --anno-type  
568 row_diff_brwt --greedy -i {metagraph_index_name}.dbg -o {metagraph_rowdiff_brtw_annotation}
```

```
569  
570 metagraph relax_brwt --relax-arity 32 -o {metagraph_rowdiff_relaxed_brtw_annotation}  
571 {metagraph_rowdiff_brtw_annotation}
```

572
573 where *{metagraph_rowdiff_relaxed_brtw_annotation}* represents the final annotation file in the row-diff relaxed
574 Multi-BRWT format.

575 The generated subindexes were subsequently compressed using the the low-level compressor XZ, using the
576 following command:

```
577  
578 xz -9 -T1 -k {uncompressed_subindex}
```

579
580 In order to calculate the overall collection size for each index, before and after compression, individual subindex
581 sizes were calculated as follows:

```
582  
583 stat -c "%s %n" {XZ-compressed/uncompressed_subindex}
```

584

585 **COBS**

586 Phylogenetically compressed COBS indexes of the 661k collection and its HQ subset are publicly available on
587 Zenodo (**Tab. S1**). Each COBS subindex of the 661k-HQ collection was downloaded from the corresponding
588 Zenodo records using Phylign (command “*make download_cobs*”), while COBS subindexes of the entire 661k
589 collection were downloaded modifying the specific https address of the Zenodo record in Phylign code and then
590 downloaded using the same command.

591 **Evaluation of search time**

592 Evaluations of search time were performed on the GenOuest cluster, using a single node with 30 CPUs and 200
593 GB of memory available for computation. Searches were performed using the EBI plasmid collection (n=2,826)
594 (**Tab. S1**) as queries. Plasmid sequences were downloaded using their corresponding plasmid IDs as follows:

```
595  
596 curl -o {plasmidID}.fasta "https://www.ebi.ac.uk/ena/browser/api/fasta/{plasmidID}"
```

597
598 and stored in a single Multi-FASTA file.

599 For each index type, we created a dedicated Snakemake workflow to perform *k*-mer matching of the EBI
600 plasmids across the 661k collection and its HQ subset, searching for all matches carrying at least 40% of query
601 *k*-mers via commands below. XZ-compressed indexes were first decompressed as follows:

602
603 `xz --decompress -k {XZ_compressed_index}`

604
605 Uncompressed indexes were queried individually. Decompression and search time was measured separately
606 for each batch via GNU Time via Galitime (**Tab. S1**). The specific commands used for k -mer matching are
607 reported below.

608 **COBS**

609 `cobs query -i {cobs_subindex} -f {EBI_plasmids} -t 0.4 -T 4 --load-complete > {output_matches}`
610

611 **Fulgor and Fulgor-m**

612 `fulgor pseudoalign -i {fulgor/meta-fulgor_subindex} -q {EBI_plasmids} -o {output_matches} --`
613 `threshold 0.4 -t 4`

614
615 where `{meta-fulgor/fulgor_subindex}` could be either a fulgor or a meta-fulgor index.

617 **Themisto-h and Themisto-r**

618 `themisto pseudoalign -i {themisto_subindex} -q {EBI_plasmids} -o {output_matches} --threshold 0.4 --`
619 `include-unknown-kmers -t 4`

620
621 where `{themisto_subindex}` could be either a Themisto-hybrid or a Themisto-roaring bitmap index, while “`--`
622 `include-unknown-kmers`” specifies that similarities should be computed considering all query k -mers.

624 **Metagraph**

625
626 `metagraph query -i {metagraph_subindex}.dbg -a {metagraph_subindex}.annodbg --count-labels --`
627 `discovery-fraction 0.4 -p 4 --fwd-and-reverse {EBI_plasmids}`

628
629 where `{metagraph_subindex}.dbg` and `{metagraph_subindex}.annodbg` represented the graph and annotation of
630 each Metagraph subindex, respectively.

633 **Phylign-Fulgor implementation and comparisons**

634 **Modified Fulgor**

635 To enable compatibility with the subindex paradigm and to integrate it in the Phylign workflow, we modified
636 the Fulgor source code (release 2.0.0, commit 12da1e9). This modified version was called “modified-Fulgor”
637 (<https://github.com/Franci-B/modified-Fulgor>) (**Tab. S1**).

- 638 • **Adjusting k -mer matching to account for all query k -mers.** Fulgor implements the so-called
639 *threshold-union* regime of k -mer matching to identify *all-matches-above-threshold*, where matches are
640 returned if they share a minimum proportion of k -mers with the query. For each match, this proportion
641 is computed by considering only positive k -mers, by default. Therefore, Fulgor’s code was modified to
642 compute this proportion considering the total number of query k -mers (modified-Fulgor, commit
643 7b040c4).
- 644 • **Returning the number of shared k -mers in the output.** Fulgor’s code was modified in order to return
645 the output in a TSV format, providing query and match IDs, along with their corresponding number of
646 shared k -mers (modified-Fulgor commit 7b040c4). Matches are output in descending order based on
647 their number of shared k -mers (modified-Fulgor commit cd2c2fa).
- 648 • **Formatting output in a COBS-like format.** To integrate Fulgor into the Phylign workflow and enable
649 the identification of best-matching genomes across different batches, the output file structure was
650 modified to follow a COBS-like format (modified-Fulgor, commit 4a6c366). To allow users to select
651 this output format, the “`--cobs`” flag was added among the input parameters (modified-Fulgor, commit
652 7b040c4).

- 653 • **Returning all matches sharing at least one k -mer.** Modifications were introduced to Fulgor's code to
654 set to 1 the minimum number of k -mers required to report matches when the similarity threshold is set
655 to 0 (modified-Fulgor, commit 669536e).
656
657

658 **Phylign-Fulgor**

659 To replace COBS with Fulgor, the Phylign repository (<https://github.com/karel-brinda/phylign>, commit
660 e1a2115) was forked to create Phylign-Fulgor (<https://github.com/Francii-B/Phylign-Fulgor>). Modified-Fulgor
661 v. 2.1.0 was integrated into Phylign as an external submodule (Phylign-Fulgor, commit 04962db).

662 Phylign-Fulgor searches were implemented for both the 661k collections (main branch) and for the ATB
663 collection (*ATB* branch). Meta-Fulgor indexes for the ATB collection were generated following the same
664 procedure as described for the 661k collection. Integration of Fulgor required the following changes in Phylign:

- 665 • **Reference collection download command.** The Fulgor-m uncompressed indexes of the 661k and
666 ATB collections were uploaded to Zenodo (**Tab. S1**). To enable downloading of Fulgor-m batches of
667 the 661k collection, download links were replaced with those of the Fulgor-m indexes both in the main
668 (Phylign-Fulgor, commit a3f915e) and in the *ATB* branch (Phylign-Fulgor, commit 7dbb010).
- 669 • **Removal of index decompression step.** The index decompression step was eliminated as
670 uncompressed Fulgor-m indexes were integrated in Phylign (Phylign-Fulgor, commit df2d7c4).
- 671 • **Perform k -mer matching using Fulgor.** COBS k -mer matching command was replaced with Fulgor's
672 corresponding command, adjusting prefixes and suffixes in filenames for referring to Fulgor indexes
673 (Phylign-Fulgor, commit df2d7c4).

674 **Sensitivity calibration**

675 To calibrate sensitivity of Phylign-Fulgor with respect to Lexicmap and Phylign-COBS, k -mer matching was
676 performed to identify the thresholds that produced comparable numbers of matched genomes when searching
677 the pCUVET18-1784 plasmid from *Serratia nevei* strain CUVET18-1784 (GenBank accession: CP115019.1;
678 length: 52,830 bp). Specifically, Phylign-Fulgor was calibrated with respect to the number of genomes matched
679 by Phylign-COBS across the 661k-HQ collection, benchmarking thresholds ranging from 5% to 40%, and with
680 respect to the genomes aligned by LexicMap across the ATB-HQ collection, benchmarking thresholds ranging
681 from 0% to 40%.

682 **Comparison with the original Phylign (Phylign-COBS)**

683 To evaluate the impact of replacing COBS with Fulgor in the Phylign, both Phylign-Fulgor and Phylign-COBS
684 were benchmarked on a MacBook Pro 18,3 (Apple M1 Pro, 10 CPU cores, 16 GB RAM). Searches were
685 performed using the EBI plasmid collection against both the 661k-HQ collection, using the following
686 parameters: `kmer_threshold=0.4, nb_best_hits=1000, minimap_preset=asm20`.

687 Sensitivity assessments were performed by searching the pCUVET18-1784 plasmid across the 661k-HQ
688 collection. Specifically, Phylign-COBS was benchmarked setting a threshold of 40% while Phylign-Fulgor was
689 tested using both its uncalibrated (40%) and calibrated thresholds (0.007% and 18%). For both Phylign-Fulgor
690 and Phylign-COBS, the following parameters were set: `nb_best_hits: 2000000, minimap_preset: asm20,`
691 `threads: 20, max_ram_gb: 80`. These experiments were executed on the GenOuest cluster, using a single node
692 with an SSD disk, 20 CPUs, and 80 GB of RAM available.

693 **Comparisons with LexicMap**

694 LexicMap index [10] for the ATB-HQ collection was built on the GenOuest cluster using the following
695 command:

```
696 lexicmap index -S -X {HQ_ATB_genome_list} -O {lexicmap_index} -b 15000 --log {log_file}
```

699 Searches with LexicMap and Phylign-Fulgor were performed both on the GenOuest cluster (single node with an
700 SSD disk, 20 CPUs, and 300 GB of RAM) and on the laptop described in the previous section. The pCUVET18-
701 1784 plasmid, previously employed in LexicMap benchmarking experiments [10], and 100 plasmids from EBI
702 database (**Tab. S1**), were searched across the ATB-HQ collection using the following LexicMap command:

```
703 lexicmap search -j 20 -d {lexicmap_index} {query} -o {output_alignments} --debug
```

706 Search time was measured using Galitime (**Tab. S1**).

707 Phylign-Fulgor searches were performed setting its calibrated threshold with respect to LexicMap (0.007%)
708 and the following parameters: *nb_best_hits*: 2000000; *minimap_preset*: asm20; *threads*: 20; *max_ram_gb*: 300.
709

710 SUPPLEMENTARY NOTES

711 Note S1. Biological question translation into a k -mer based problem

712 Given a biological question of interest, we assume it can be answered using k -mers, and assume that it can be
713 formalized using a k -mer based pseudo-algorithm involving a step of k -mer matching. Its specific formalization
714 may call for different *matching strategies*, i.e. a specific combination of these three elements: *matching mode*,
715 characteristics of the reference collection, and query type.

716 As an example, we describe two questions recently discussed in literature (**Fig 1a**), highlighting the role of
717 k -mer matching and the specific matching strategies employed for addressing the problem.

718 The first example was provided by Bradley and colleagues [17], who used k -mers to characterize plasmid
719 host-range in reference collection. Their bioinformatic method involved an initial step of k -mer matching
720 finalized to the detection of putative conjugative systems by searching for exact matches between the alleles of
721 two specific types of conjugative genes, namely MOB and T4SS, and the sequences in the reference collection.
722 Post-processing analysis followed k -mer matching to determine the presence of these putative conjugative
723 systems, defined by the presence of at least one allele for both types of genes.

724 The second example was provided by Brinda and colleagues [8], who developed the Genomic Neighbour
725 Typing (GNT) technique for predicting the phenotypic antimicrobial resistance profile of bacterial isolates. In
726 this case, the step of k -mer matching is performed by ProPhyle whose role is to identify the “nearest neighbors”
727 of the bacterial isolate of interest in a reference collection. The matching strategy used in GNT was to search the
728 set of best matching sequences for each sequencing long-reads of the bacterial isolate. Post-processing analysis
729 followed k -mer matching to merge the results of each reads and to predict the isolate’s AMR phenotypic profile
730 based on the phenotypic profile of its nearest neighbors.

731 Note S2. k -mer index selection

732 Selection of the most suitable k -mer index for a given application should be driven first by its compatibility with
733 both the subindex paradigm and the defined matching strategy, and then by its compression and searching
734 performances. Ideally, these evaluations should allow users to identify the best k -mer index that best minimizes
735 both the search time of k -mer matching and the use of computational resources. If a single optimal solution for
736 both qualities is not available among the candidates, a Pareto optimization analysis should be performed to
737 decide whether to choose the best trade-off index or a combination of multiple k -mer index types. To validate
738 the final choice (either a single or multiple indexes combined together), the chosen k -mer index(es) has to be
739 integrated in the overall bioinformatic workflow and tested as well.

740 Note S3. Supported matching modes

741 Both Metagraph and COBS support all matching modes through their threshold-based methods (defined as *exact*
742 *k -mer matching* and *approximate pattern matching* by their respective authors), where thresholds are computed
743 as a fraction of matched k -mers over the total number of query k -mers. Specifically for the *all-best-matches* and
744 *all-best-matches-above-threshold* modes, an additional post-processing step is required to filter the output
745 matches and retain only the best matches. This is possible, as both indexes report the list of query-match pairs
746 together with the corresponding number of shared k -mers.

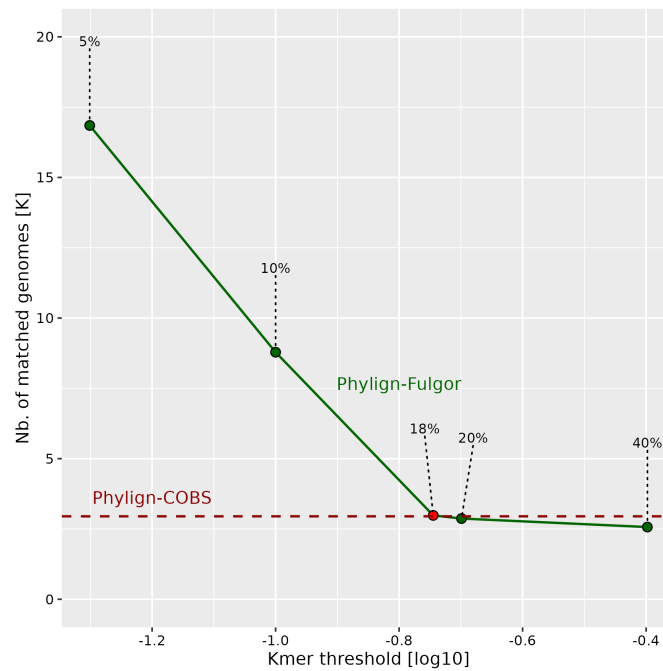
747 In contrast, Themisto and Fulgor partially support the *all-matches-above-threshold* mode, as they both
748 implement threshold-based searches, defined as *threshold* and *threshold-union* methods by their respective
749 authors. In both tools, similarity can be computed in one of two ways: either considering the full set of query k -
750 mers, or only the positive subset, namely those query k -mers that have been found at least once in the index. To
751 ensure compatibility with the *subindex paradigm*, it is essential to choose the former option (i.e. using all query
752 k -mers). Themisto supports this via the parameter “*--include-unknown-kmers*”, while Fulgor requires further
753 minor modifications in its code to enable this option. However, this mode is only partially supported because the
754 output does not report how many k -mers the query shares with each match, making it impossible to rank
755 matches based on their “similarity” with the query.

756 The remaining matching modes are not supported by Themisto and Fulgor. *Exact-matches* cannot be
757 identified because 1) they both implement the so-called *intersection* method that ignores k -mers that are entirely
758 absent from the index, and 2) neither return query-match similarity information in the output. Therefore, it is not
759 possible to distinguish true exact matches (i.e., those that contain all query k -mers) from matches lacking the k -
760 mers absent from the index. Similarly, *All-matches-above-threshold* and *all-matches* cannot be identified by
761 Themisto and Fulgor, as 1) neither tool implements a dedicated method for these matching modes, and 2) they
762 do not output query-match similarity values, which prevents these modes from being derived from either
763 intersection- or threshold-based search.

764 **SUPPLEMENTARY FIGURES**

765

766



767

768 **Fig. S1. k -mer threshold calibration for Phylogen-Fulgor's *matches-above-threshold* using a**
769 **single plasmid and 661k-HQ**

770 The graph compares the number of matched genomes as a function of the required proportion of matching k -
771 mers in the *all-matches-above-threshold* mode, in comparison with Phylogen-COBS. The red point ($t=18\%$)
772 corresponds to the threshold yielding a similar number of matches to Phylogen-COBS.

773

774 **SUPPLEMENTARY TABLES**

775 **Tab. S1. Overview of the software and data used throughout the paper**

776

Type	Name	Details	URL
Data collections	EBI plasmid collection	FASTA sequences	https://www.ebi.ac.uk/genomes/plasmid
	100 EBI plasmids	FASTA sequences	https://github.com/Franciai-B/optimized-kmer-search-supplement
	pCUVET18-1784 plasmid	FASTA sequences	https://www.ncbi.nlm.nih.gov/nuccore/CP115019.1
	661k	metadata	https://dx.doi.org/10.6084/m9.figshare.16437939
		tar.xz assemblies	https://zenodo.org/records/4602622
		Phylogenetically compressed COBS indexes	https://zenodo.org/records/7313926 https://zenodo.org/records/6849657 https://zenodo.org/records/7315499
	661k-HQ	Phylogenetically compressed COBS indexes	https://doi.org/10.5281/zenodo.6849657 https://doi.org/10.5281/zenodo.6845083
		Phylogenetically compressed Fulgor-m indexes	https://zenodo.org/records/14002973 https://zenodo.org/records/14002975 https://zenodo.org/records/14006705 https://zenodo.org/records/14006707
	ATB	tar.xz assemblies	https://osf.io/xv7q9/
	ATB-HQ	Phylogenetically compressed Fulgor-m indexes	https://zenodo.org/records/15994164 https://zenodo.org/records/15994228 https://zenodo.org/records/15994270 https://zenodo.org/records/15994318 https://zenodo.org/records/15994445 https://zenodo.org/records/15994505 https://zenodo.org/records/15994553 https://zenodo.org/records/15994624
Software	COBS	v.0.3.1	https://github.com/bingmann/cobs
	Metagraph	v.0.3.6	https://github.com/ratschlab/metagraph
	Themisto	v.3.2.2	https://github.com/algbio/themisto
	Fulgor	v.2.0.0	https://github.com/jermp/fulgor
	modified-Fulgor	v.2.1.0	https://github.com/Franciai-B/modified-Fulgor
	Phylign-COBS	commit e1a2115	https://github.com/karel-brinda/Phylign
	Phylign-Fulgor	<ul style="list-style-type: none"> ● main commit 72fe2e4 ● ATB commit e635c0c 	https://github.com/Franciai-B/Phylign-Fulgor
	LexicMap	v. 0.7.0	https://github.com/shenwei356/LexicMap
	Galitime	v. 0.1.3	https://github.com/karel-brinda/galitime

777

778 **Tab. S2. Detailed statistics for BLAST-like alignments by Phylign-COBS and Phylign-Fulgor**

779 Search time, memory requirements and outputs of the pipelines for searching all EBI plasmids (n=2,826) across
780 the 661k-HQ collection. Both Phylign-COBS and Phylign-Fulgor were run on a MacBook Pro 18,3, setting the
781 following parameters: *kmer_thres*: 0.4; *nb_best_hits*: 1000; *minimap_preset*: *asm20*; *threads*: 4.
782

Method		Computational time [h]		Matching statistics			
		Real time	CPU time	#distinct query-genome pairs	#matched queries	#target genomes	#target batches
Phylign-Fulgor	match	2.13	8.43	3,732,698	1,740	330,683	295
	align	36.54	13.61	716,790	1,738	193,568	291
Phylign-COBS	match	3.44	34.97	9,163,123	1,873	367,565	298
	align	13.20	14.92	838,830	1,871	205,231	296

783

784 **Tab. S3. Alignment comparisons between Phylign-Fulgor and Phylign-COBS**

785 Search time, memory usage, and the number of hits obtained when searching the pCUVET18-1784 plasmid
 786 sequence (length: 52,830 bp) across the 661k-HQ dataset (n. genomes=639,981, n. batches=305) are reported
 787 for both tools. Phylign-Fulgor and Phylign-COBS were run on the GenOuest cluster
 788 (<https://www.genouest.org>), using a single node with an SSD disk, 20 CPUs, and 80 GB of RAM available. The
 789 following parameters were used for Phylign-Fulgor: *nb_best_hits*: 2,000,000; *minimap_preset*: *asm20*; *threads*:
 790 20; *max_ram_gb*: 80.

Method		Matching statistics		Server (20 cores)		
		#target genomes	#alignments	Real time [s]	CPU time [s]	RAM [GB]
<i>Single plasmid (pCUVET18-1784)</i>						
Phylign-Fulgor	match (t=0.007%)	136,777	-	115	230	0.5
	match+align(t=0.007%)	114,775	414,865	3,467	57,788	3.5
	match+align(t=18%)	2,982	-	70	145	0.2
	match+align(t=18%)	2,982	26,079	617	5,803	2.8
	match (t=40%)	2,568	-	128	176	0.6
	match+align(t=40%)	2,568	22,956	642	5,245	2.4
Phylign-COBS	match (t=40%)	2,949	-	583	6,969	61.2
	match+align(t=40%)	2,949	25,870	1,086	12,581	61.2

791