

1 **MetaCerberus: distributed highly parallelized scalable HMM-based**
2 **implementation for robust functional annotation across the tree of**
3 **life**

4 Jose L. Figueroa III¹⁻², Eliza Dhungel¹, Cory R. Brouwer¹, Richard Allen White III^{1-2*}

5 ¹North Carolina Research Campus (NCRC), Department of Bioinformatics and
6 Genomics, The University of North Carolina at Charlotte, 150 Research Campus Drive,
7 Kannapolis, NC 28081, USA

8 ²Computational Intelligence to Predict Health and Environmental Risks (CIPHER),
9 Department of Bioinformatics and Genomics, The University of North Carolina at
10 Charlotte, 9201 University City Boulevard, Charlotte, NC 28223, USA

11

12

13 *To whom correspondence should be addressed

14

15 **Keywords:** Microbiomes, Microbes, Viruses, Bacteriophage (phage), Metaomics,
16 Annotation, HMMs, HMMER, GAGE, Pathview

17

18

19

20

21

22

23

24

25 **Abstract**

26 **Summary:** MetaCerberus is an exclusive HMM/HMMER-based tool that is massively
27 parallel, on low memory, and provides rapid scalable annotation for functional gene
28 inference across genomes to metacommunities. It provides robust enumeration of
29 functional genes and pathways across many current public databases including KEGG
30 (KO), COGs, CAZy, FOAM, and viral specific databases (i.e., VOGs and PHROGs). In a
31 direct comparison, MetaCerberus was twice as fast as EggNOG-Mapper, and produced
32 better annotation of viruses, phages, and archaeal viruses than DRAM, PROKKA, or
33 InterProScan. MetaCerberus annotates more KOs across domains when compared to
34 DRAM, with a 186x smaller database and a third less memory. MetaCerberus is fully
35 integrated with differential statistical tools (i.e., DESeq2 and edgeR), pathway
36 enrichment (GAGE R), and Pathview R for quantitative elucidation of metabolic
37 pathways. MetaCerberus implements the key to unlocking the biosphere across the tree
38 of life at scale.

39 **Availability and implementation:** MetaCerberus is written in Python and distributed under
40 a BSD-3 license. The source code of MetaCerberus is freely available at
41 <https://github.com/raw-lab/metacerberus>. Written in python 3 for both Linux and Mac OS
42 X. MetaCerberus can also be easily installed using `mamba create -n metacerberus -c`
43 `bioconda -c conda-forge metacerberus`

44
45 **Contact:** Richard Allen White III, UNC Charlotte, rwhit101@charlotte.edu

46 **Supplementary information:** Supplementary data are available online.

47

48 **Introduction**

49 Annotation is a fundamental step in functional gene inference, which is required
50 by many disciplines in biology. Massively parallel sequencing (MPS) has reached the
51 terabyte scale with Illumina NovaSeq X producing 16 Tb per run and Oxford nanopore
52 promethION 7 Tb per run (**1-2**). Due to this increase in MPS, the number of reference
53 microbial genomes and metagenomes has increased by orders of magnitude. Genome
54 Taxonomy Database (GTDB) now includes 402,709 (08-RS214, April 28th, 2023)
55 genomes, and the Short Read Archive (SRA) has >4.5 million listed biosample
56 metagenomes (**3-4**). Cellular metagenome-assembled genomes (MAGs) and their viral
57 counterpart vMAGs (viral MAGs) have also rapidly populated public databases through
58 reconstruction from shotgun metagenomics (**5-7**). Functional gene annotation is
59 required for metabolic reconstruction, functional and structural gene differential analysis,
60 inference of pathway regulation, presence/absence of toxin genes (e.g., botulinum toxin
61 A), novel gene cluster discovery (e.g., antibiotic discovery), and viral detection. Due to
62 this Terabyte scale, the annotation step will be the most prolonged, requiring more CPU
63 time, memory, and resources to finish before obtaining biological insight. Reference
64 databases have also been nearing the Terabyte scale, taking days to download and
65 format, requiring massive allocations of disk space. Thus scalable, highly parallel, low
66 memory, and rapid annotation tools are critical to the future of 'omics analysis.

67 Functional annotation requires two main steps: 1) gene calling followed by 2)
68 gene assignment via external reference databases. Multiple approaches have been
69 applied for gene calling and gene assignment, including homology and ontology-based
70 methods. Gene calling finds protein-coding open reading frames (pORFs) alongside

71 ribosomal RNAs, transfer RNAs, and other RNAs. Various tools exist for pORF calling,
72 including Prodigal, FragGeneScanRs, GetOrf, and GeneMark (**8-11**). Gene assignment
73 of pORFs to external databases often uses homology-based tools such as BLAST (**12**),
74 MMseq2 (**13**), and/or DIAMOND (**14**) against databases such as RefSeq (NCBI
75 Reference Sequence Database) (**15**), UniProt (Universal Protein Resource) (**16**), or
76 KEGG (Kyoto Encyclopedia of Genes and Genomes) (**17**). Common tools include
77 PROKKA, DRAM (Distilled and Refined Annotation of Metabolism), InterProScan
78 (INTEgrative PROtein signature database), EggNOG-Mapper (evolutionary genealogy
79 of genes: Non-supervised Orthologous Groups), and MicrobeAnnotator (**18-22**).
80 Ontology-based approaches are generally superior to homology-based methods (**21**).
81 EggNOG-Mapper and InterProScan utilize homology-based (i.e., Diamond and BLAST)
82 and Hidden Markov Models (HMMs) based ontology approaches via HMMER (**23**) using
83 either KEGG (**17**), eggNOG (**24**), InterPro (**25**), or Pfam (protein family) databases (**26**).
84 HMMs provide greater sensitivity to elucidate and discover relationships between query
85 and database based on ontology and are protein domain-centric (**21, 27**).

86 Viruses and the candidate phyla radiation (CPR) have remained challenging to
87 functionally annotate due to the divergent nature of their proteins (**28-29**). DRAM has a
88 specific version (i.e., DRAM-v) to annotate viruses, including the detection of viral
89 auxiliary metabolic genes (vAMGs) (**19**). MicrobeAnnotator and DRAM have attempted
90 to close the gap in CPR functional annotation. While no specific annotation tool or gene
91 database exists for CPR, they are found amongst GTDB and other public repositories
92 (**30**). Various databases such as VOGs (Virus Orthologous Groups), pVOGs
93 (Prokaryotic Virus Orthologous Groups), IMG/VR (Integrated Microbial Genome/Virus),

94 INPHARED (INfrastructure for a PHAge REference Database), and PHROGs
95 (Prokaryotic Virus Remote Homologous Groups database) have been introduced to
96 improve annotation viruses from isolates and vMAGs (**31-35**). Still, CPR and viruses
97 remain a significant challenge for functional annotation.

98 Many tools are available for functional annotation from genomes to
99 metagenomes; however, gaps remain between: 1) resource utilization (e.g., memory
100 use), 2) large database size, and 3) parallel processing, and simultaneously providing
101 robust rapid annotation at scale. Further development of tools for CPR and viral
102 functional annotation are a general community need. We present MetaCerberus, an
103 ontology-based HMM tool that provides scalable, highly parallel, low memory usage,
104 and rapid annotation for genomes to metacommunities across the tree of life.

105 **Implementation**

106 *Framework and coding base*

107 MetaCerberus is written entirely in Python (version 3) as a wrapper for various
108 other tools described below. Similar to our other software MerCat2 for massively parallel
109 processing (MPP), it utilizes a byte chunking algorithm 1 ('Chunker') to split files for
110 MPP for further utilization in RAY, a massive open-source parallel computing framework
111 to scale Python applications and workflows (**36**). Using RAY's scalable parallelization
112 within MetaCerberus allows utilization across multiple nodes with ease. To avoid large
113 RAM consumption, we implemented the greedy algorithm for tab-separated merging
114 and incremental PCA plot limiter from MerCat2 (**36**). MetaCerberus utilizes
115 HMM/HMMER exclusively without homology-based tools (e.g., BLAST)

116

117 *Databases for MetaCerberus*

118 MetaCerberus enables functional gene assignment across multiple databases,
119 including: 1) KOfams (KEGG protein families) to obtain KEGG KOs (KEGG Ontology)
120 (version 11-Jul-2023, <https://www.genome.jp/ftp/db/kofam/>), 2) FOAM (Functional Ontology
121 Assignments for Metagenomes), 3) COG (Clusters of Orthologous Genes) (version
122 2020, <https://ftp.ncbi.nih.gov/pub/COG/COG2020/data/>), and 4) dbCAN (DataBase for
123 automated Carbohydrate-active enzyme ANnotation) for CAZy (Carbohydrate-Active
124 enZymes Database) (version 11, <https://bcb.unl.edu/dbCAN2/download/>) (**37-41**). For viral
125 annotation, MetaCerberus enlists VOG (version 219, <https://vogdb.org/download>), pVOG
126 (version Sep2016, <https://ftp.ncbi.nlm.nih.gov/pub/kristensen/pVOGs/downloads.html#>), and
127 PHROG (version 4, <https://phrogs.lmge.uca.fr/>) databases. FOAM ontology is obtained
128 from KOfam KOs, and then computed via a reference table to avoid redundancy.
129 Similarly, the dbCAN database is used to obtain CAZy ontology via a reference table.
130 COGs and PHROGs are currently not formatted as HMMs within their public
131 repositories. We converted them into protein family-specific HMMs (e.g., COG1 ->
132 COG1.hmm) using MAFFT (version 7.273-woe) via local alignments with maximum
133 iterations of 1000 (**42**). We compared databases of six other tools to MetaCerberus,
134 including DRAM, PROKKA, InterProScan, MicrobeAnnotator, and EggNOG-Mapper
135 (**Table 1**). Currently, only MetaCerberus provides functional annotation and support to
136 FOAM, pVOG, and PHROG databases (**Table 1**). EggNOG-Mapper and MetaCerberus
137 are the only tools we compared that supports the COG database (**Table 1**). All tools
138 compared in this study obtain the enzyme commission numbers (EC) numbers (**Table**
139 **1**).

140 *Modes for running MetaCerberus*

141 MetaCerberus has three basic modes: 1) quality control (QC) for raw reads, 2)
142 formatting/gene prediction, and 3) annotation (**Fig 1**). MetaCerberus can use three
143 different input files: 1) raw read data from any sequencing platform (Illumina, PacBio, or
144 Oxford Nanopore), 2) assembled contigs, as MAGs, vMAGs, isolate genomes, or a
145 collection of contigs, 3) amino acid fasta (.faa), previously called pORFs (**Fig 1**). We
146 offer customization, including running databases all together, individually or specifying
147 select databases. For example, if a user wants to run a prokaryotic or eukaryotic-
148 specific KOfam, or an individual database alone such as dbCAN, both are easily
149 customized within MetaCerberus. In future versions, we will provide viral and phage-
150 specific KO modules to run individually. In QC mode, raw reads are quality controlled
151 via fastqc (version v0.12.1) prior and post trim (**43**). Raw reads are then trimmed via
152 data type; if the data is Illumina or PacBio, fastp (version 0.23.4) is called, otherwise it
153 assumes the data is Oxford nanopore then Porechop (version v0.2.4) is utilized (**43-45**,
154 **Fig 1**). Post quality-control trimmed reads are converted to fasta without quality (**Fig 1**).
155 If Illumina reads are utilized, an optional bbmap (version 39.01) step to remove the
156 phiX174 genome is available. Phage phiX174 is a common contaminant within the
157 Illumina platform as it is their library spike-in control (**46-47**). We highly recommend this
158 removal if viral analysis is conducted, as it would provide false positives to ssDNA
159 microviruses within the sample.

160 In the formatting and gene prediction mode, contigs and genomes are checked
161 for N repeats. These N repeats are removed by default. We impute contig/genome
162 statistics (e.g., N50, N90, max contig) via our custom module Metaome Stats. Contigs

163 are converted to pORFs via Prodigal or FragGeneScanRs as specified by user
164 preference (**48, Fig 1**). Scaffold annotation is not recommended due to N's providing
165 ambiguous annotation. Both callers can be used via our --super option, and we
166 recommend using FragGeneScanRs for samples rich in eukaryotes as it performed
167 better in our hands than Prodigal (unpublished data). HMMER searches against the
168 above databases via user specified bitscore and e-values or our minimum defaults (i.e.,
169 bitscore = 25, e-value = 1×10^{-9}).

170 There are six general rules followed by MetaCerberus for functional annotation.
171 Rule 1 is the *score pre-filtering module* for pORFs thresholds: each pORF match to an
172 HMM is recorded by default or user-selected e-value/bit scores per database
173 independently, across all databases, or per user specification of the selected database.
174 Rule 2 is imputed for *non-overlapping dual domain module* pORF threshold: if two HMM
175 hits are non-overlapping from the same database, both are counted as long as they are
176 the within the default or user selected e-value/bit scores. Rule 3 is computed as the
177 *winner take all module* for overlapping pORFs: if two HMM hits are overlapping (>10
178 amino acids) from the same database the lowest resulting e-value and highest bit score
179 wins. Rule 4 is *similar match independent accession module* for a single pORF: if both
180 hits within the same database have equal values for both e-value and bit score but are
181 different accessions from the same database (e.g., KO1 and KO3) then both are
182 reported. Rule 5 is the *whole count incomplete exclusion module* filter only allows
183 whole discrete integer counting. Rule 6 the *dual independent overlapping domain*
184 *module* for convergent binary domain pORFs. If two domains within a pORF are
185 overlapping <10 amino acids (e.g, COG1 and COG4) then both domains are counts and

186 reported due to the dual domain issue within a single pORF. If a function hits multiple
187 pathways within an accession, both are counted, in pathway roll-up, as many proteins
188 function in multiple pathways.

189

190 *Statistics and visualization*

191 MetaCerberus, as previously mentioned, provides genome and contig statistics
192 via MetaOme stats; it also offers seamless integration into automatic differential
193 statistics, visualizations, pathway enrichment, and pathway integration viewing. DESeq2
194 and edgeR negative binomial distribution differential statistic tools are available to users
195 by selection (default is DESeq2) (**49-50**). The outputs from DESeq2, edgeR, or both are
196 automatically enriched for pathway analysis in GAGE (Generally Applicable Gene-set
197 Enrichment for Pathway Analysis) R (**51**). GAGE outputs are loaded into Pathview R to
198 visualize differential pathways across user-specified experiments (**52**). These outputs
199 include differential KEGG heatmaps from Pathview, volcano plots, and gene level
200 heatmaps (**Fig S1**).

201 A sample dashboard visualization is provided for all data input types (e.g., reads,
202 contigs and/or genomes) with a number of pORF called, MetaOme stats (i.e., genome
203 statistics, N50, N90, etc., for genomes/contigs only), PCA with sample sets of >3, and
204 the number of annotated hits for all databases or user select specifications (**Fig S1**).

205

206 *Across tool comparisons*

207 Tools compared across MetaCerberus (version 1.1) include DRAM (version
208 1.4.6), InterProScan (version 5.60-92.0), EggNog-Mapper (version 2.1.8),

209 MicrobeAnnotator (version 2.0.5), and PROKKA (version 1.1). All comparisons were
210 completed on a Dual 8-Core Intel Xeon E5-2667 CPU @ 3.2GHz (16 cores) using
211 128GB RAM. MPP testing of MetaCerberus was completed on five nodes of a Dual 18-
212 Core Intel Xeon Gold 6154 CPU @ 3.00GHz (36 cores/node). All genomes used in our
213 study are available at <https://osf.io/3uz2j/>. For further testing of MetaCerberus, we used
214 five distinct genospecies, *Rhizobium leguminosarum*, against five distinct
215 *Exiguobacterium* spp. available at
216 https://github.com/raw-lab/MetaCerberus/tree/main/data/rhizobium_test (**Table S1**).
217 Viruses from permafrost that were used in the DRAM paper
218 (<https://www.ncbi.nlm.nih.gov/nuccore/QGNH000000000>) were compared directly to
219 MetaCerberus and DRAM (**19**).

220

221 *Data availability*

222 Sequence files, genome files, and supplemental data are available at
223 <https://osf.io/3uz2j/>. Databases are also freely available at <https://osf.io/3uz2j/>. All code is
224 available at www.github.com/raw-lab/metacerberus.

225

226 *Contributing to MetaCerberus and Fungene*

227 MetaCerberus is a community resource as is recently acquired FunGene
228 (<http://fungene.cme.msu.edu/>). We welcome contributions of other experts expanding
229 annotation of all domains of life (viruses, phages, bacteria, archaea, eukaryotes).
230 Please send us an issue on our MetaCerberus GitHub.
231 (www.github.com/raw-lab/metacerberus/issue); we will fully annotate your genome, add

232 suggested pathways/metabolisms of interest, make custom HMMs to be added to
233 MetaCerberus and FunGene.

234

235 **Results**

236 *Database size and download time*

237 Formatting and downloading are required steps in functional annotation and both
238 depend on database size. Substantial databases take up large amounts of costly disk
239 space and require expensive computers with large amounts of costly RAM for analysis.
240 MetaCerberus database size is 3.8 GB, with a download time of ~4 mins, and database
241 format time is zero because they are pre-formatted already for the user (**Table 2**).
242 DRAM database download requires 710 GB of disc space, and requires ~3 days to
243 download completely (**Table 2**). According to the DRAM readme, KEGG Genes and
244 UniRef90 need ~500 GB of disc space and ~512 GB of RAM to process the complete
245 database (**19**, <https://github.com/WrightonLabCSU/DRAM>). This database size difference is
246 due UniRef90 updates since their original release in 2020. DRAM can run with more
247 processors within a single node but is not set up for multi-node like MetaCerberus. The
248 InterProScan database is 14 GB, which took ~2.45 h to install (**Table 2**). PROKKA had
249 the smallest database at 636 MB and had the fastest install of ~3 ½ minutes (**Table 2**).
250 MicrobeAnnotator requires at least ~237 GB for its full version and ~0.65 GB for its light
251 version (**Table 2, 22**).

252

253

254

255 *Computational resource comparison*

256 We compared MetaCerberus to DRAM, InterProScan, and PROKKA for the time
257 used per genome, RAM utilization, and disk space used across 100 randomly selected
258 bacterial genomes within GTDB (**Table S1, Fig 2**). Generally, PROKKA had the highest
259 processing speed per genome (~48 sec median, **Fig 2**). InterProScan had the slowest
260 at ~21 min per genome median time (**Fig 2**). DRAM was ~5 mins faster per genome
261 than MetaCerberus (i.e., 10 mins vs 15 mins) (**Fig 2**). MetaCerberus and PROKKA had
262 the lowest RAM, followed by InterProScan (**Fig 2**). DRAM using default parameters had
263 the highest RAM observed (**Fig 2**). DRAM had the lowest disc space due to the deletion
264 of files post-finalization, with PROKKA having the most disc space (**Fig 2**). EggNOG-
265 mapper using HMMs was initially compared to MetaCerberus; however, further testing
266 was not completed due to the high run time failure rate. On average, EggNOG-mapper
267 failed to finish annotation 32% of the time using 148 randomly selected GTDB
268 bacterial/archaeal genomes used by other tools (**Table S3**). Approximately 16% of the
269 genomes failed after running for two days; another 16% could not annotate even when
270 other tools, including MetaCerberus, had no issue (**Table S3**). The average annotation
271 time with EggNOG-mapper v2 was ~53 min, with a median of ~30 min (**Table S3**). It
272 was the slowest tool tested and thus removed from further comparisons.
273 MicrobeAnnotator didn't functionally install the database correctly. We could not
274 successfully use the code; thus, it was removed from further comparisons..

275

276

277

278 *Automatic statistical and pathway analysis*

279 MetaCerberus provides automatic differential statistics, pathway gene
280 enrichments, and KEGG map-based heatmaps in Pathview R for data exploration, data
281 mining, and hypothesis generation. As a test, we compared five distinct genospecies,
282 *Rhizobium leguminosarum*, against five distinct *Exiguobacterium* spp. using
283 MetaCerberus using both DESeq2 and edgeR (**Table S3**). These genomes were
284 selected as a comparison due to differential pathways within the comparison genomes.
285 *Rhizobium* are symbiotic nitrogen fixers containing both nitrogenase for nitrogen fixation
286 and nodulation genes for symbiotic nodule formation within legume roots (**53**). The
287 *Exiguobacterium* spp. Have a bright orange colony morphology color from
288 biosynthetically made carotenoids; it is hypothesized that the carotenoid
289 diaponeurosporene-4-oic acid from the C₃₀ carotenoid biosynthesis pathway is what
290 provides the distinctive orange color (**54**). Chemical studies have suggested other
291 carotenoids are present within *Exiguobacterium* spp. that contribute to the orange
292 colony color (**55**). MetaCerberus found differential pathway assignments using DESeq2
293 and Pathview for carotenoid biosynthesis, ABC transporters, and phosphotransferase
294 system (including nitrogen regulation) (**Fig S2-4**). edgeR found an additional pathway in
295 benzoate degradation that wasn't found in DESeq2 (**Fig S5**).

296

297 *Annotation comparisons*

298 PROKKA, DRAM, and MetaCerberus all use Prodigal for pORF calling.
299 MetaCerberus also provides an extra pORF caller FragGeneScanRs. InterProScan
300 uses the EMBOSS getorf pORF caller, which in all cases had lower pORFs than

301 Prodigal regardless of the genome kingdom type (e.g., bacteria, archaea, CPR, phage,
302 archaeal virus or eukaryotic virus) (**Fig S6**). Generally, PROKKA, DRAM, and
303 MetaCerberus had similar pORF calling numbers; however, DRAM did call more pORF
304 from eukaryotic viruses (**Fig S6**).

305 Furthermore, we compared MetaCerberus to DRAM, InterProScan, and
306 PROKKA for whether a pORF was annotated, listed as hypothetical, or unknown (no
307 annotation). We randomly selected 100 unique bacterial genomes from GTDB, 100
308 unique archaea genomes from GTDB, 100 unique phage genomes from INPHARED,
309 100 unique eukaryotic viral genomes from RefSeq, 78 CPR genomes, and 82 archaeal
310 viral genomes for these annotation tests (**Table S1**). MetaCerberus, DRAM, and
311 InterProScan protein modes had similar annotation results of ~78-83% for bacteria, with
312 InterProScan being the highest at 83% (**Fig 3**). InterProScan using nucleotide mode
313 had the lowest annotation amount across all kingdoms (**Fig 3-4**). PROKKA had ~50% of
314 the pORFs as annotated and hypothetical for bacteria and ~60% hypothetical for
315 archaea (**Fig 3**). CPR annotation InterProScan had the highest at 70%, followed by
316 DRAM at 66%, then MetaCerberus at 61% (**Fig 3**). MetaCerberus and PROKKA had
317 fewer unknowns than DRAM for bacteria, archaea, and CPR genomes (**Fig 3**).
318 PROKKA annotated very few CPR pORFs, with the majority >60% being hypothetical
319 proteins (**Fig 3**). DRAM generally doesn't find many hypothetical proteins or lists them
320 as unknown across domains of the tree of life.

321 MetaCerberus performs better for viruses, phages, and archaeal viruses (**Fig 4**).
322 MetaCerberus annotates more per genome >63 % phages, >65 % viruses, and >41%
323 archaeal viruses based on median values (**Fig 4**). MetaCerberus outperforms across all

324 viruses (e.g., phages, eukaryotic viruses, and archaeal viruses) providing more
325 annotations, fewer hypotheticals, and fewer unknowns compared to DRAM,
326 InterProScan, and PROKKA (**Fig 4**). MetaCerberus annotates more KOs from KOfams
327 than DRAM across all domains (**Fig 5**). PROKKA and InterProScan don't provide KOs;
328 therefore, we couldn't compare KOs found across domains to MetaCerberus.

329 To better compare to DRAM-v, the only other tool exclusively for viruses and
330 phages, we analyzed a virome containing 1,907 viral populations (VPs) obtained from
331 Swedish permafrost. Based on time, MetaCerberus took 99 mins to complete the
332 annotation compared to 141.75 mins for DRAM (**Fig 6**). When MetaCerberus is utilized
333 at it's full potential with RAY it only takes 12.5 mins for the same dataset (**Fig 6**). RAM
334 was significant less with MetaCerberus vs. DRAM-v, both in MPP and non-MPP mode,
335 with <500 Mb of RAM compared to 18.7 Gb with DRAM-v (**Fig 6**). MetaCerberus had
336 more annotations than DRAM-v for the Swedish permafrost virome across shared
337 databases (i.e., KO, CAZy, and VOG) (**Fig S7**).

338

339 **Discussion**

340 MetaCerberus provides a low memory, robust, scalable, and rapid annotation
341 across the tree of life, exclusively using HMMs/HMMER. HMMER is a powerful tool to
342 find pORFs that may be missed by standard homology-based tools due to its protein
343 domain centric and supervised machine-learning nature. It's rarely used alone due to
344 the speed and time required to finish annotation. MetaCerberus has provided a solution
345 to this scaling issue using RAY and algorithms needed from MerCat2. EggNOG-Mapper
346 v2 is the only other tool that exclusively provides HMMs/HMMER-based annotation

347 alone. MetaCerberus runs twice as fast on a single node than EggNOG-Mapper v2
348 without RAY. MetaCerberus with RAY is three times as fast as EggNOG-Mapper v2
349 (data not shown) on a database that is 1/3 smaller.

350 Generally, MetaCerberus performs better for viral and phage annotation when
351 directly compared to DRAM-v. DRAM finds more pORFs than MetaCerberus (**Fig S6**)
352 due to it using the -meta option in Prodigal for viruses; whereas, MetaCerberus uses a
353 standard for complete genomes in this case but still can annotate viral genomes better
354 than DRAM on a much smaller database. Viruses, archaeal viruses, and phages are a
355 grand challenge to unlock the 'unknown' and 'hypothetical' functions within their
356 genomes.

357 As data scales, computational time, memory, and waiting for results will take
358 longer. Scalable tools like MetaCerberus are needed as we approach Petabyte levels of
359 sequencing. MetaCerberus provides a further community resource to annotate the
360 unknowns of our biosphere. Lastly, MetaCerberus provides a robust tool kit to annotate
361 the entire tree of life at scale.

362

363 **Funding:** R.A. White III and Jose Figueroa are supported by the UNC Charlotte
364 Department Bioinformatics and Genomics start-up package from the North Carolina
365 Research Campus in Kannapolis, NC, and by NSF ABI Development award 1565030.
366 Cory Brouwer and Eliza Dhungel were also supported by NSF ABI Development award
367 1565030.

368

369 **Acknowledgments:** We also acknowledge the University Research Computing and the
370 College of Computing and Informatics for computational and logistical support. We must
371 further acknowledge J Peter W Young at the University of York for our discussions on
372 the nature of *Rhizobium* genomes, especially genospecies. For his insightful grammar
373 edits and proofreading of multiple drafts, we acknowledge Bryan W. Fulghum.

374

375 **Conflicts of Interest:** The authors declare no conflicts of interest. RAW is the CEO of
376 RAW Molecular Systems (RAW), LLC, but no financial, IP, or others from RAW LLC
377 were used or contributed to the study.

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403 References

404

405 1. Illumina throughput specs (date accessed July 17th, 2023).

406 <https://www.illumina.com/systems/sequencing-platforms/novaseq-x-plus.html>

407

408 2. Oxford throughput specs (date accessed July 17th, 2023).

409 [https://nanoporetech.com/about-us/news/highest-throughput-yet-promethion-breaks-7-terabase-](https://nanoporetech.com/about-us/news/highest-throughput-yet-promethion-breaks-7-terabase-mark)

410 mark

411

412 3. Genome Taxonomy Database (GTDB) release statistics (date accessed July 17th,

413 2023).

414 <https://gtdb.ecogenomic.org/>

415

416 4. Short Read Archive Biosample Metagenomes (date accessed July 17th, 2023).

417 <https://www.ncbi.nlm.nih.gov/sra/?term=metagenomes>

418

419 5. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK,

420 Schulz F, Jarett J, Rivers AR, Eloë-Fadrosh EA, Tringe SG, Ivanova NN, Copeland A,

421 Clum A, Becraft ED, Malmstrom RR, Birren B, Podar M, Bork P, Weinstock GM, Garrity

422 GM, Dodsworth JA, Yooseph S, Sutton G, Glockner FO, Gilbert JA, Nelson WC, Hallam

423 SJ, Jungbluth SP, Etema TJG, Tighe S, Konstantinidis KT, Liu WT, Baker BJ, Rattei T,

424 Eisen JA, Hedlund B, McMahon KD, Fierer N, Knight R, Finn R, Cochrane G, Karsch-

425 Mizrachi I, Tyson GW, Rinke C, Genome Standards C, Lapidus A, Meyer F, Yilmaz P,

426 Parks DH, Eren AM, Schriml L, Banfield JF, Hugenholtz P, Woyke T. Minimum

427 information about a single amplified genome (MISAG) and a metagenome-assembled

428 genome (MIMAG) of bacteria and archaea. *Nat Biotechnol.* 2017;35(8):725–731.

429

430 6. Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, Kuhn

431 JH, Lavigne R, Brister JR, Varsani A, Amid C, Aziz RK, Bordenstein SR, Bork P,

432 Breitbart M, Cochrane GR, Daly RA, Desnues C, Duhaime MB, Emerson JB, Enault F,

433 Fuhrman JA, Hingamp P, Hugenholtz P, Hurwitz BL, Ivanova NN, Labonté JM, Lee KB,

434 Malmstrom RR, Martinez-Garcia M, Mizrachi IK, Ogata H, Páez-Espino D, Petit MA,

435 Putonti C, Rattei T, Reyes A, Rodriguez-Valera F, Rosario K, Schriml L, Schulz F,

436 Steward GF, Sullivan MB, Sunagawa S, Suttle CA, Temperton B, Tringe SG, Thurber

437 RV, Webster NS, Whiteson KL, Wilhelm SW, Wommack KE, Woyke T, Wrighton KC,

438 Yilmaz P, Yoshida T, Young MJ, Yutin N, Allen LZ, Kyrpides NC, Eloë-Fadrosh EA.

439 Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nat Biotechnol.*

440 2019 37(1):29-37.

441

442 7. Kieft K, Adams A, Salamzade R, Kalan L, Anantharaman K. vRhyme enables binning

443 of viral genomes from metagenomes. *Nucleic Acids Res.* 2022;50(14):e83.

444

445 8. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:

446 prokaryotic gene recognition and translation initiation site identification. *BMC*

447 *Bioinformatics.* 2010;11:119.

448

- 449 9. Van der Jeugt F, Dawyndt P, Mesuere B. FragGeneScanRs: faster gene prediction
450 for short reads. *BMC Bioinformatics*. 2022 3(1):198.
451
- 452 10. Emboss getorf (date accessed July 17th, 2023).
453 <https://emboss.sourceforge.net/apps/cvs/emboss/apps/getorf.html>
454
- 455 11. Besemer J, Borodovsky M. GeneMark: web software for gene finding in
456 prokaryotes, eukaryotes and viruses. *Nucleic Acids Res*. 2005 33:W451-4.
457
- 458 12. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden
459 TL. BLAST+: architecture and applications. *BMC Bioinform*. 2009;10:421.
460
- 461 13. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for
462 the analysis of massive data sets. *Nat Biotechnol*. 2017 35(11):1026-1028.
463
- 464 14. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using
465 DIAMOND. *Nat Methods*. 2015, 12(1):59–60.
466
- 467 15. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B,
468 Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretin A, Bao Y, Blinkova O,
469 Brover V, Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T,
470 Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P,
471 McGarvey KM, Murphy MR, O'Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD,
472 Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan
473 AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A, Tatusova T, DiCuccio M, Kitts P,
474 Murphy TD, Pruitt KD. Reference sequence (RefSeq) database at NCBI: current status,
475 taxonomic expansion, and functional annotation. *Nucleic Acids Res*. 2016 44(D1):D733-
476 45.
477
- 478 16. UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic*
479 *Acids Res*. 2023 51(D1):D523-D531.
480
- 481 17. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new
482 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res*.
483 2017;45:D353–61.
484
- 485 18. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. 2014.
486 30(14):2068-9.
487
- 488 19. Shaffer M, Borton MA, McGivern BB, Zayed AA, La Rosa SL, Solden LM, Liu P,
489 Narrowe AB, Rodríguez-Ramos J, Bolduc B, Gazitúa MC, Daly RA, Smith GJ, Vik DR,
490 Pope PB, Sullivan MB, Roux S, Wrighton KC. DRAM for distilling microbial metabolism
491 to automate the curation of microbiome function. *Nucleic Acids Res*. 2020 48(16):8883-
492 8900.
493
- 494 20. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J,
495 Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong

- 496 SY, Lopez R, Hunter S. InterProScan 5: genome-scale protein function classification.
497 Bioinformatics. 2014 30(9):1236-40.
498
- 499 21. Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J.
500 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain
501 Prediction at the Metagenomic Scale. *Mol Biol Evol.* 2021 38(12):5825-5829.
502
- 503 22. Ruiz-Perez CA, Conrad RE, Konstantinidis KT. MicrobeAnnotator: a user-friendly,
504 comprehensive functional annotation pipeline for microbial genomes. *BMC*
505 *Bioinformatics.* 2021 22(1):11.
506
- 507 23. Eddy SR. Accelerated Profile HMM Searches. *PLoS Comput Biol.* 2011
508 7(10):e1002195.
509
- 510 24. Hernández-Plaza A, Szklarczyk D, Botas J, Cantalapiedra CP, Giner-Lamia J,
511 Mende DR, Kirsch R, Rattei T, Letunic I, Jensen LJ, Bork P, von Mering C, Huerta-
512 Cepas J. eggNOG 6.0: enabling comparative genomics across 12 535 organisms.
513 *Nucleic Acids Res.* 2023 51(D1):D389-D394.
514
- 515 25. Paysan-Lafosse T, Blum M, Chuguransky S, Grego T, Pinto BL, Salazar GA,
516 Bileschi ML, Bork P, Bridge A, Colwell L, Gough J, Haft DH, Letunić I, Marchler-Bauer
517 A, Mi H, Natale DA, Orengo CA, Pandurangan AP, Rivoire C, Sigrist CJA, Sillitoe I,
518 Thanki N, Thomas PD, Tosatto SCE, Wu CH, Bateman A. InterPro in 2022. *Nucleic*
519 *Acids Res.* 2023 January 6th;51(D1):D418-D427.
520
- 521 26. Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL,
522 Tosatto SCE, Paladin L, Raj S, Richardson LJ, Finn RD, Bateman A. Pfam: The protein
523 families database in 2021. *Nucleic Acids Res.* 2021 49(D1):D412-D419
524
- 525 27. Steinegger M, Meier M, Mirdita M, Vöhringer H, Haunsberger SJ, Söding J. HH-
526 suite3 for fast remote homology detection and deep protein annotation. *BMC*
527 *Bioinformatics.* 2019 20(1):473.
528
- 529 28. Fremin BJ, Bhatt AS, Kyrpides NC; Global Phage Small Open Reading Frame (GP-
530 SmORF) Consortium. Thousands of small, novel genes are predicted in global phage
531 genomes. *Cell Rep.* 2022 39(12):110984.
532
- 533 29. Jaffe AL, Castelle CJ, Matheus Carnevali PB, Gribaldo S, Banfield JF. The rise of
534 diversity in metabolic platforms across the Candidate Phyla Radiation. *BMC Biol.* 2020
535 18(1):69.
536
- 537 30. Parks DH, Chuvochina M, Rinke C, Mussig AJ, Chaumeil PA, Hugenholtz P. GTDB:
538 an ongoing census of bacterial and archaeal diversity through a phylogenetically
539 consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Res.*
540 2022 50(D1):D785-D794.
541

- 542 31. Virus Orthologous Groups (VOG) database
543 <https://vogdb.org/>
544
- 545 32. Graziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups
546 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic*
547 *Acids Res.* 2017 January 4th;45(D1):D491-D498. doi: 10.1093/nar/gkw975. Epub 2016
548 October 26th. PMID: 27789703; PMCID: PMC5210652.
549
- 550 33. Camargo AP, Nayfach S, Chen IA, Palaniappan K, Ratner A, Chu K, Ritter SJ,
551 Reddy TBK, Mukherjee S, Schulz F, Call L, Neches RY, Woyke T, Ivanova NN, Eloie-
552 Fadrosh EA, Kyrpides NC, Roux S. IMG/VR v4: an expanded database of uncultivated
553 virus genomes within a framework of extensive functional, taxonomic, and ecological
554 metadata. *Nucleic Acids Res.* 2023 51(D1):D733-D743.
555
- 556 34. Cook R, Brown N, Redgwell T, Rihtman B, Barnes M, Clokie M, Stekel DJ, Hobman
557 J, Jones MA, Millard A. INfrastructure for a PHAge REference Database: identification
558 of large-scale biases in the current collection of cultured phage genomes. *Phage (New*
559 *Rochelle).* 2021 2(4):214-223.
560
- 561 35. Terzian P, Olo Ndela E, Galiez C, Lossouarn J, Pérez Bucio RE, Mom R, Toussaint
562 A, Petit MA, Enault F. PHROG: families of prokaryotic virus proteins clustered using
563 remote homology. *NAR Genom Bioinform.* 2021 3(3):lqab067.
564
- 565 36. Figueroa III JL, Panyala A, Colby S, Friesen ML, Tiemann L, White III RA. MerCat2:
566 a versatile k-mer counter and diversity estimator for database-independent property
567 analysis obtained from omics data. *bioRxiv*, 2022.
568
- 569 37. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H.
570 KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive score
571 threshold. *Bioinformatics.* 2020 36(7):2251-2252.
572
- 573 38. Prestat E, David MM, Hultman J, Taş N, Lamendella R, Dvornik J, Mackelprang R,
574 Myrold DD, Jumpponen A, Tringe SG, Holman E, Mavromatis K, Jansson JK. FOAM
575 (Functional Ontology Assignments for Metagenomes): a Hidden Markov Model (HMM)
576 database with environmental focus. *Nucleic Acids Res.* 2014 42(19):e145.
577
- 578 39. Galperin MY, Wolf YI, Makarova KS, Vera Alvarez R, Landsman D, Koonin EV.
579 COG database update: focus on microbial diversity, model organisms, and widespread
580 pathogens. *Nucleic Acids Res.* 2021 49(D1):D274-D281.
581
- 582 40. Yin Y, Mao X, Yang J, Chen X, Mao F, Xu Y. dbCAN: a web resource for automated
583 carbohydrate-active enzyme annotation. *Nucleic Acids Res.* 2012;40:W445–51.
584
- 585 41. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. The
586 carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res.*
587 2014;42:D490–5.
588

- 589 42. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
590 improvements in performance and usability. *Mol Biol Evol.* 2013 30(4):772-80.
591
- 592 43. FASTQC
593 <https://github.com/s-andrews/FastQC>
594
- 595 44. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
596 *Bioinformatics.* 2018 34(17):i884-i890.
597
- 598 45. Porechop
599 <https://github.com/rrwick/Porechop>
600
- 601 46. Moustafa A, Xie C, Kirkness E, Biggs W, Wong E, Turpaz Y, Bloom K, Delwart E,
602 Nelson KE, Venter JC, Telenti A. The blood DNA virome in 8,000 humans. *PLoS*
603 *Pathog.* 2017 13(3):e1006292.
604
- 605 47. Bushnell, Brian. BBMap: A Fast, Accurate, Splice-Aware Aligner. United States.
606
- 607 48. MetaOme Stats
608 https://github.com/raw-lab/metaome_stats
609
- 610 49. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion
611 for RNA-seq data with DESeq2. *Genome Biol.* 2014 15(12):550.
612
- 613 50. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for
614 differential expression analysis of digital gene expression data. *Bioinformatics.* 2010
615 26(1):139-40.
616
- 617 51. Luo W, Friedman MS, Shedden K, Hankenson KD, Woolf PJ. GAGE: generally
618 applicable gene set enrichment for pathway analysis. *BMC Bioinformatics.* 2009 10:161.
619
- 620 52. Luo W, Brouwer C. Pathview: an R/Bioconductor package for pathway-based data
621 integration and visualization. *Bioinformatics.* 2013 29(14):1830-1.
622
- 623 53. Young JPW, Moeskjær S, Afonin A, Rahi P, Maluk M, James EK, Cavassim MIA,
624 Rashid MH, Aserse AA, Perry BJ, Wang ET, Velázquez E, Andronov EE, Tampakaki A,
625 Flores Félix JD, Rivas González R, Youseif SH, Lepetit M, Boivin S, Jorin B, Kenicer
626 GJ, Peix Á, Hynes MF, Ramírez-Bahena MH, Gulati A, Tian CF. Defining the *Rhizobium*
627 *leguminosarum* Species Complex. *Genes (Basel).* 2021 12(1):111.
628
- 629 54. White RA 3rd, Soles SA, Gavelis G, Gosselin E, Slater GF, Lim DSS, Leander B,
630 Suttle CA. The Complete Genome and Physiological Analysis of the Eurythermal
631 Firmicute *Exiguobacterium chiriquicha* Strain RW2 Isolated From a Freshwater
632 Microbialite, Widely Adaptable to Broad Thermal, pH, and Salinity Ranges. *Front*
633 *Microbiol.* 2019 9:3189.
634

635 55. Jinendiran S, Dahms HU, Dileep Kumar BS, Kumar Ponnusamy V, Sivakumar N.
636 Diapolycopenedioic-acid-glucosyl ester and keto-myxocoxanthin glucoside ester: Novel
637 carotenoids derived from *Exiguobacterium acetylicum* S01 and evaluation of their
638 anticancer and anti-inflammatory activities. *Bioorg Chem.* 2020 103:104149.
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680

681 **Figure and Table Legends**

682 **Table 1:** Comparing tools based on databases provided. This includes versions of other
683 databases present within the various tools compared.

684

685 **Table 2:** Database size, download, and formatting time across tools.

686

687 **Figure 1:** Flowgraph of the MetaCerberus pipeline. MetaCerberus has three input data
688 formats which include raw reads, contigs, or previously called pORFs. Quality control
689 step is mainly utilized for raw reads only. Contigs have a formatting mode were they are
690 quality controlled for N presence, followed by N removal, and also provides basic contig
691 statistics using Metaome Stats (e.g., N50, N90 etc). Gene calling currently offers
692 prodigal or FragGeneScanRs for pORF calling. Gene prediction is completed with
693 HMMER/HMM against KEGG and FOAM KOs (all by default). Users can select
694 additional databases such as CAZy, COG, PHROG, VOG, and pVOGs for viruses, or
695 selective KOfams for prokaryotes or eukaryotes. With running four or more samples it
696 provides a PCA for KEGG and FOAM KOs, a basic run metric dashboard, as well as
697 differential statistics using DESeq2/edgeR, and pathway enrichment using GAGE R
698 followed by plotting in Pathview R.

699

700 **Figure 2:** Computational resource comparison. DRAM, InterProScan, PROKKA, and
701 MetaCerberus are compared computationally for time to complete each genome
702 annotation, RAM required to complete annotation per genome, and disc space needed
703 for inputs/outputs. The 100 randomly selected bacterial genomes were from GTDB
704 (**Table S1**).

705

706 **Figure 3:** Annotation comparison across cellular domains of life (bacteria, archaea,
707 CPR). MetaCerberus was compared to DRAM, InterProScan, and PROKKA for
708 annotation across various genomes. Supplemental materials include the genomes for
709 bacteria, archaea, and CPR used in this comparison (**Table S1**).

710

711 **Figure 4:** Annotation comparison across viruses infecting differential cellular domains
712 (phage, archaeal viruses, eukaryotic viruses). MetaCerberus was compared to DRAM,
713 InterProScan, and PROKKA for annotation across various genomes. The genomes are
714 listed for phage, archaeal viruses, and eukaryotic viruses in supplemental materials
715 (**Table S1**).

716

717 **Figure 5.** DRAM vs. MetaCerberus KO annotation comparison across the domains of
718 life. DRAM and MetaCerberus utilize KOfams for KEGG KO assignment if the user
719 doesn't provide a KEGG KO database separately. The genomes for the comparison are
720 listed in supplemental materials (**Table S1**). The e-values and bitscore can vary
721 between DRAM and MetaCerberus. In this comparison, we choose the default dbCAN
722 e-value option of $<1e^{-15}$ and the default bitscore of 60 for DRAM.

723

724 **Figure 6.** DRAM vs. MetaCerberus computational resource comparison. A virome from
725 Swedish permafrost containing 1,907 VPs were compared computationally for time to

726 complete annotation, RAM required to complete annotation, and disc space needed for
727 inputs/outputs. MPP testing for MetaCerberus utilized five nodes for comparisons.

728

729 **Supplemental Materials**

730

731 **Table S1:** List of genomes used for computational comparisons. This includes
732 randomly selected GTDB genomes for archaea and bacteria domains, phage genomes,
733 archaeal viral genomes, CPR genomes, and RefSeq viral genomes.

734

735 **Figure S1:** Standard output dashboard for MetaCerberus. Outputs include a complete
736 html drawn in plotly. Also, for comparisons of >3 genomes, FOAM and KEGG-based KO
737 PCA are included.

738

739 **Figure S2:** Comparisons *Rhizobium* vs. *Exiguobacterium* genomes using MetaCerberus
740 for carotenoid pathways in Pathview. KO counts from KEGG were normalized with
741 DESeq2, enriched with GAGE, then plotted with Pathview R. Genomes are listed in
742 supplemental materials (**Table S1**).

743

744 **Figure S3:** Comparisons *Rhizobium* vs. *Exiguobacterium* genomes using MetaCerberus
745 for ABC transporters in Pathview. KO counts from KEGG were normalized with
746 DESeq2, enriched with GAGE, then plotted with Pathview R. Genomes are listed in
747 supplemental materials (**Table S1**).

748

749 **Figure S4:** Comparisons *Rhizobium* vs. *Exiguobacterium* genomes using MetaCerberus
750 for phosphotransferase system in Pathview. KO counts from KEGG were normalized
751 with DESeq2, enriched with GAGE, then plotted with Pathview R. Genomes are listed in
752 supplemental materials (**Table S1**).

753

754 **Figure S5:** Comparisons *Rhizobium* vs. *Exiguobacterium* genomes using MetaCerberus
755 for Benzoate degradation in Pathview. KO counts from KEGG was normalized with
756 edgeR, enriched with GAGE, then plotted with Pathview R. Genomes are listed in
757 supplemental materials (**Table S1**).

758

759 **Figure S6:** Comparisons of protein-coding open reading frame (pORF) calling across
760 computational tools. All tools but InterProScan use Prodigal for pORFs. InterProScan
761 uses the EMBOSS getorf tool.

762

763 **Figure S7:** Comparing DRAM-v vs. MetaCerberus for annotation across shared
764 databases (i.e., KO, VOG, CAZy). The Swedish virome was utilized for comparisons.

765

766

767

768

769

770

771

772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822

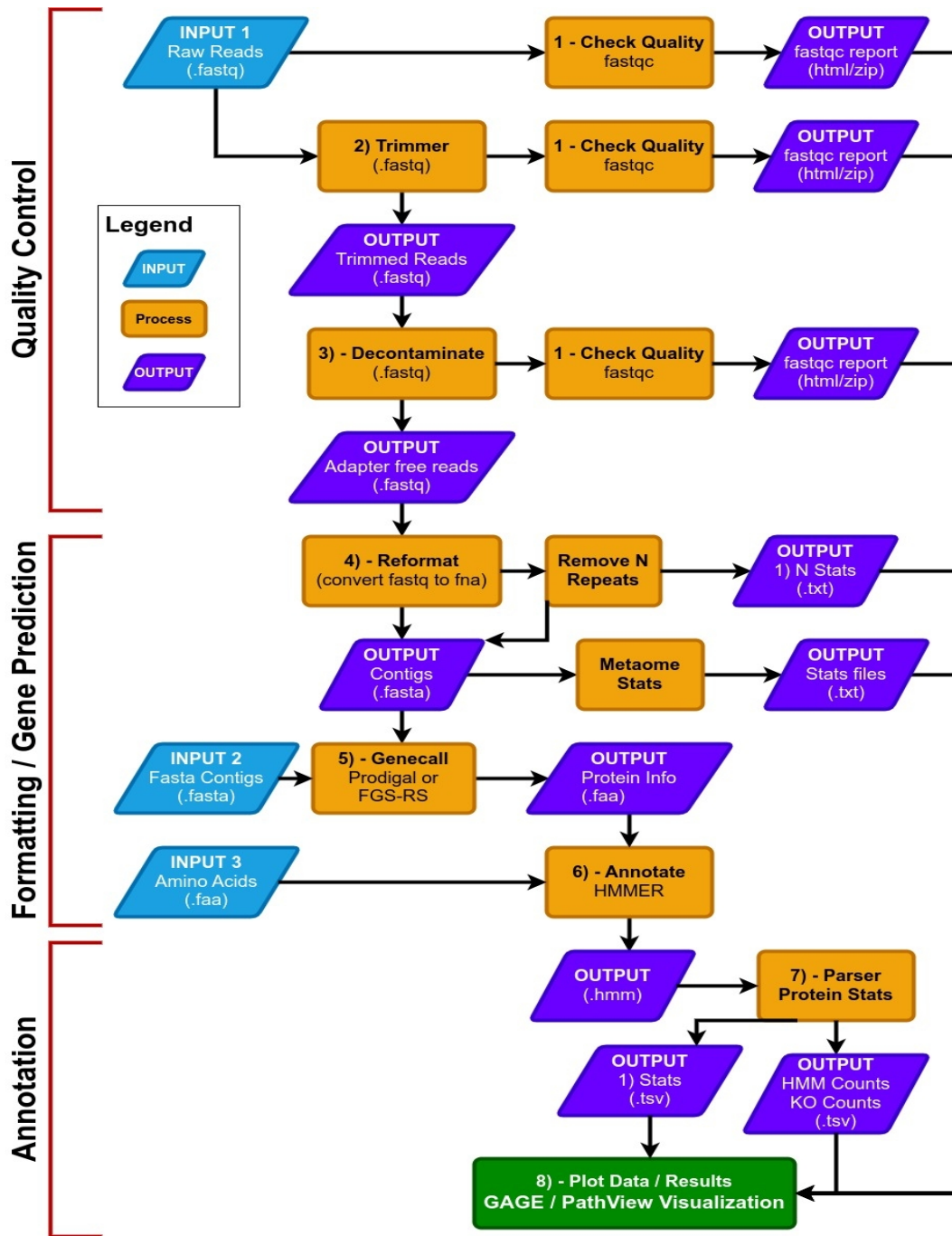


Figure 1.

823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873

GTDB-bacteria

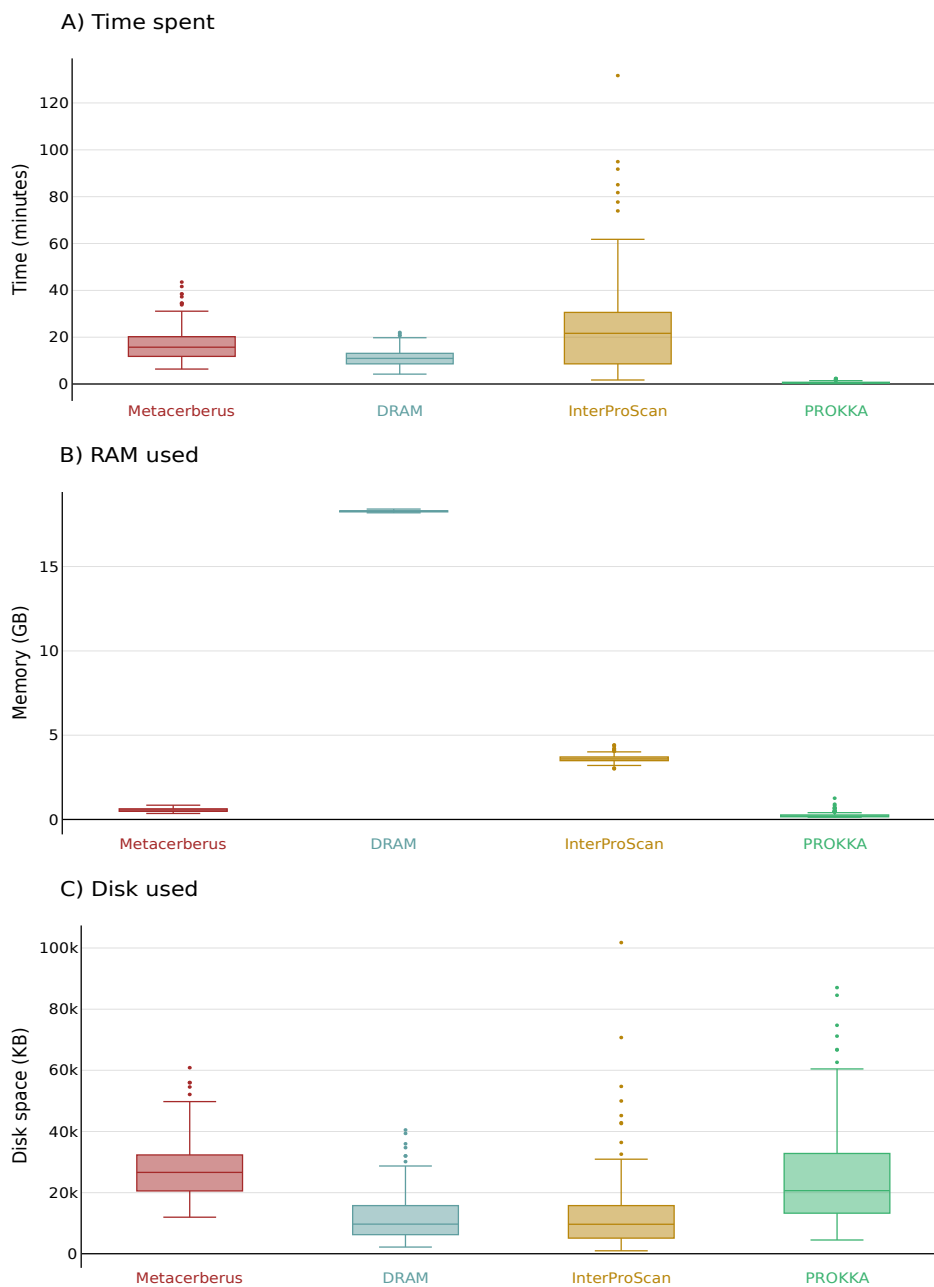
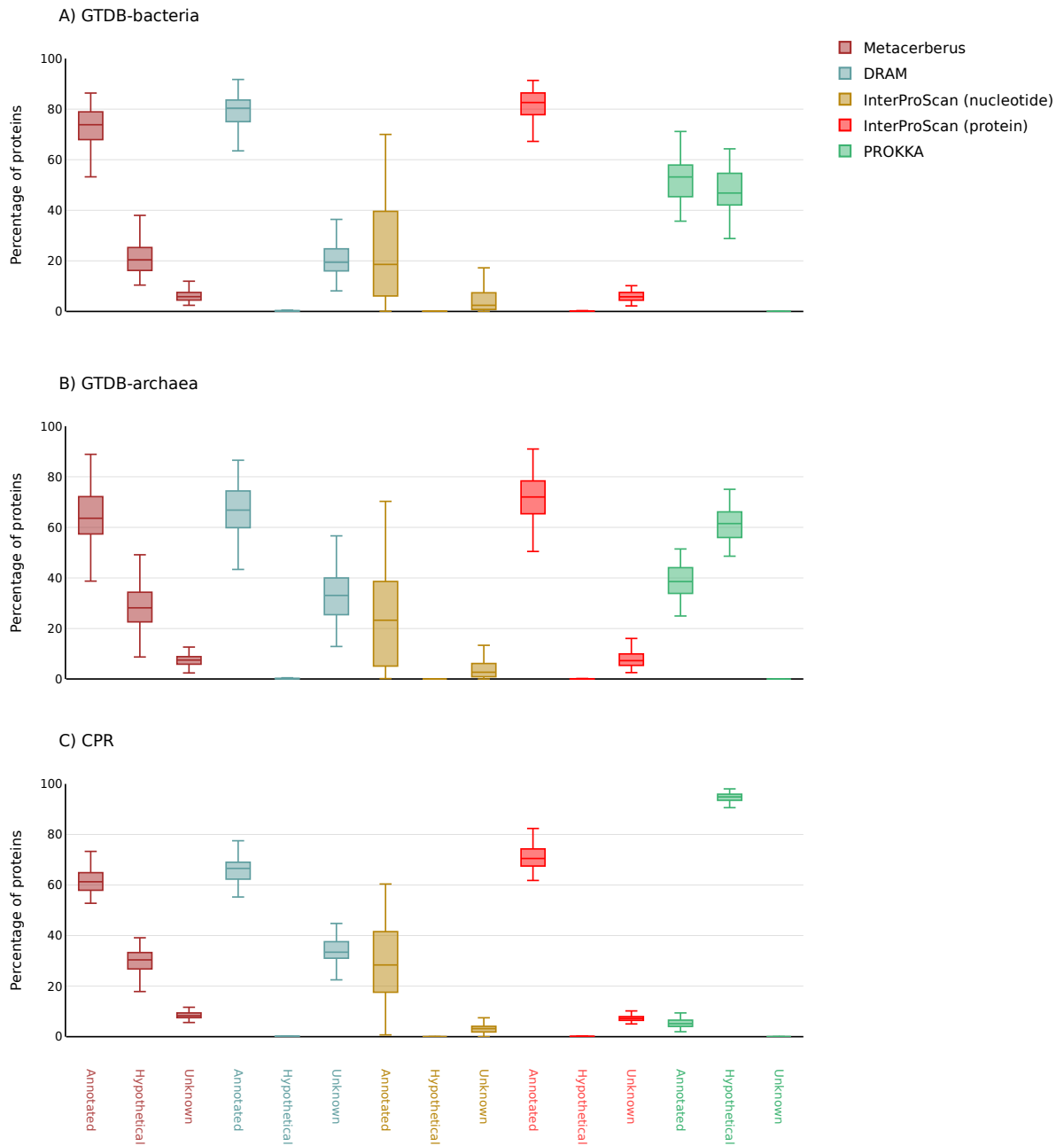
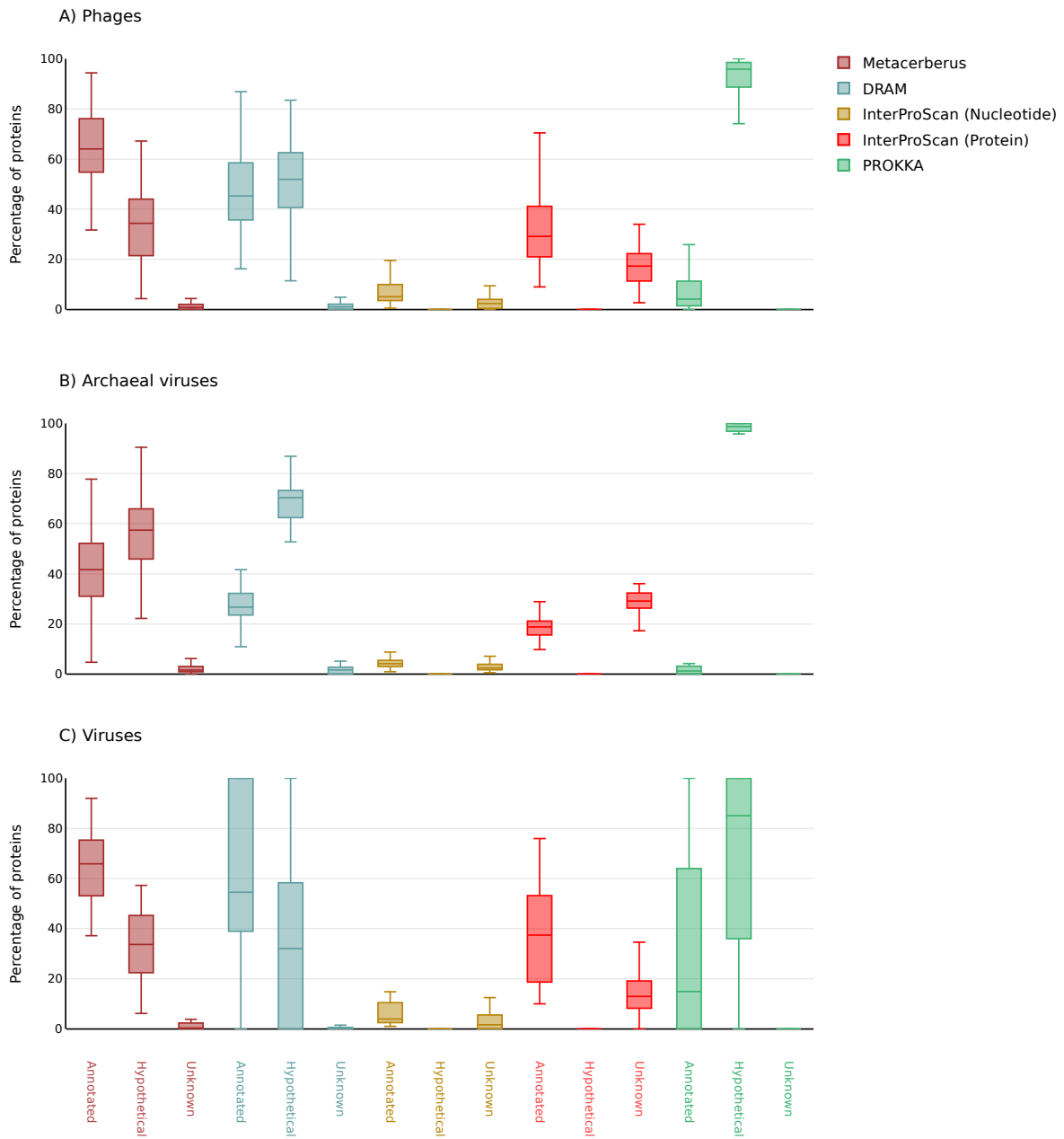


Figure 2.



875
876 **Figure 3.**
877
878
879
880
881
882
883

884



886 **Figure 4.**

887

888

889

890

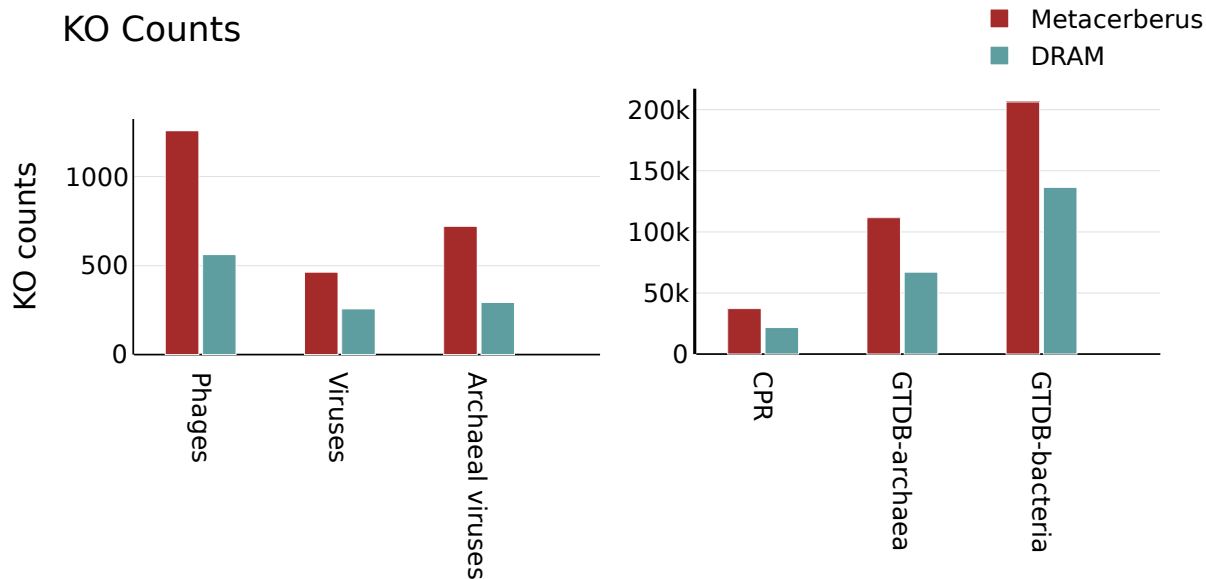
891

892

893

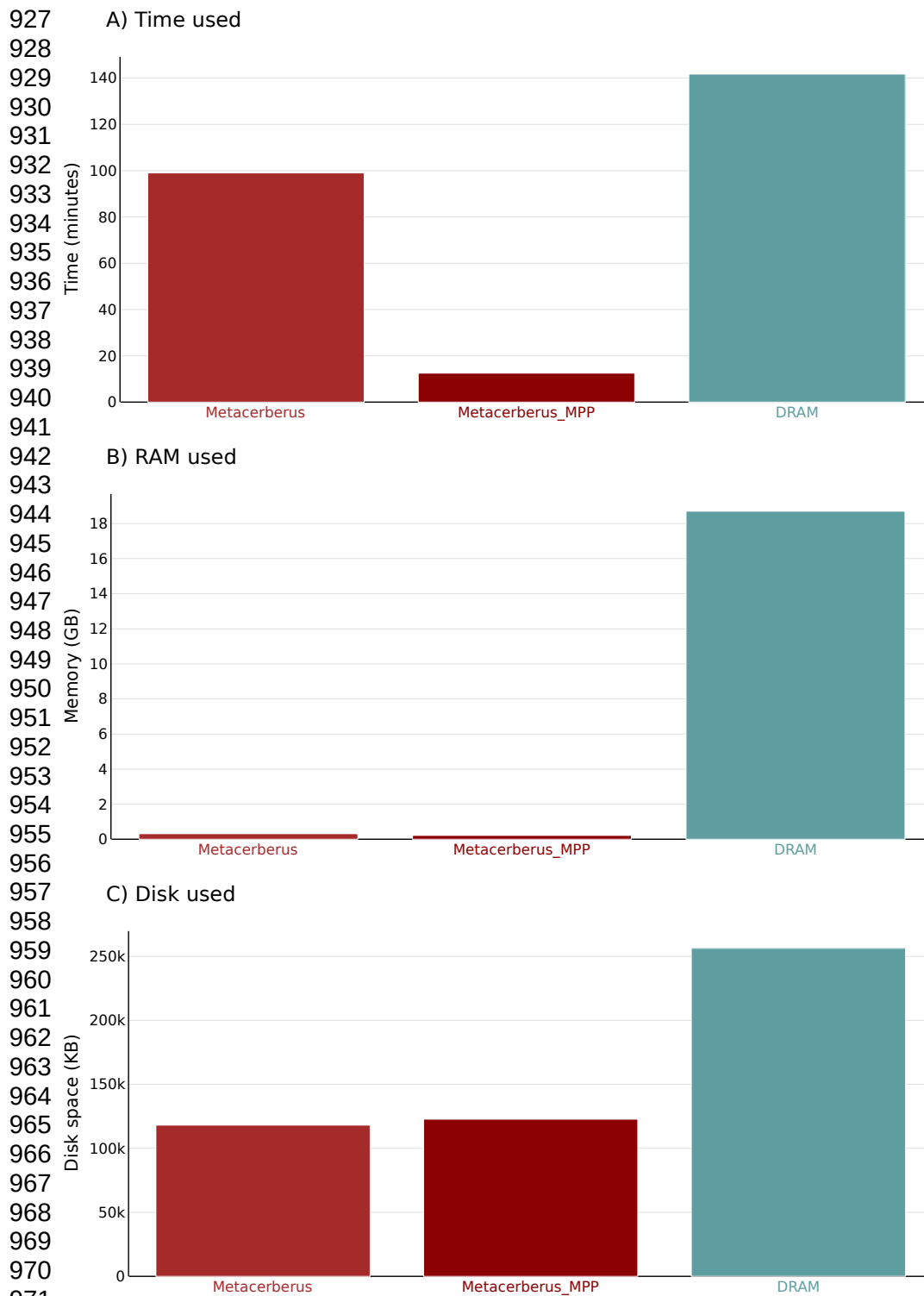
894

895



897 **Figure 5.**

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
924
925
926



972
973
974 **Figure 6.**
975
976
977

978 **Table 1**

	EC	KEGG	CAZy	COG	FOAM	VOG	pVOG	PHROG	pfam	EggNOG	InterPro
MetaCerberus	X	X	X	X	X	X	X	X			
DRAM	X	X	X			X			X		
Prokka	X								X		
InterProScan	X								X		X
MicrobeAnnotator	X	X									X
EggNOG-Mapper	X	X	X	X					X	X	

980
981

Table 2

Tool	Time	Disk	Version
DRAM	~ 3 days	~710GB	v1.4.6
InterProScan	~ 2:45:59.23	14GB	v5.60-92.0
Metacerberus	~ 0:04:14.29	3.8GB	v1.1
PROKKA	~ 0:03:28.68	607M	v1.14.6
EggNOG-Mapper	~14:33:31.74	31GB	V2.1.8
MicrobeAnnotator	>3 days	~237GB	v2.0.5

983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001